Low-Rank Approximation with $1/\epsilon^{1/3}$ Matrix-Vector Products

Ainesh Bakshi abakshi@cs.cmu.edu CMU Kenneth L. Clarkson klclarks@us.ibm.com IBM David P. Woodruff dwoodruf@cs.cmu.edu CMU

Abstract

We study iterative methods based on Krylov subspaces for low-rank approximation under any Schatten-p norm. Here, given access to a matrix ${\bf A}$ through matrix-vector products, an accuracy parameter ${\boldsymbol \epsilon}$, and a target rank k, the goal is to find a rank-k matrix ${\bf Z}$ with orthonormal columns such that $\|{\bf A}({\bf I}-{\bf Z}{\bf Z}^{\top})\|_{\mathcal{S}_p} \leq (1+\epsilon) \min_{{\bf U}^{\top}{\bf U}={\bf I}_k} \|{\bf A}({\bf I}-{\bf U}{\bf U}^{\top})\|_{\mathcal{S}_p}$, where $\|{\bf M}\|_{\mathcal{S}_p}$ denotes the ℓ_p norm of the the singular values of ${\bf M}$. For the special cases of p=2 (Frobenius norm) and $p=\infty$ (Spectral norm), Musco and Musco (NeurIPS 2015) obtained an algorithm based on Krylov methods that uses $\tilde{O}(k/\sqrt{\epsilon})$ matrix-vector products, improving on the naïve $\tilde{O}(k/\epsilon)$ dependence obtainable by the power method, where $\tilde{O}(\cdot)$ suppresses poly $(\log(dk/\epsilon))$ factors.

Our main result is an algorithm that uses only $\tilde{O}(kp^{1/6}/\epsilon^{1/3})$ matrix-vector products, and works for *all*, not necessarily constant, $p \ge 1$. For p = 2 our bound improves the previous $\tilde{O}(k/\epsilon^{1/2})$ bound to $\tilde{O}(k/\epsilon^{1/3})$. Since the Schatten-p and Schatten- ∞ norms of any matrix are the same up to a $1 + \epsilon$ factor when $p \ge (\log d)/\epsilon$, our bound recovers the result of Musco and Musco for $p = \infty$. Further, we prove a matrix-vector query lower bound of $\Omega(1/\epsilon^{1/3})$ for any fixed constant $p \ge 1$, showing that surprisingly $\tilde{\Theta}(1/\epsilon^{1/3})$ is the optimal complexity for constant k.

To obtain our results, we introduce several new techniques, including optimizing over *multiple Krylov subspaces* simultaneously, and *pinching inequalities* for partitioned operators. Our lower bound for $p \in [1,2]$ uses the *Araki-Lieb-Thirring* trace inequality, whereas for p > 2, we appeal to a *norm-compression* inequality for *aligned partitioned operators*. As our algorithms only require matrix-vector product access, they can be applied in settings where alternative techniques such as sketching cannot, e.g., to covariance matrices, Hessians defined implicitly by a neural network, and arbitrary polynomials of a matrix.

1 Introduction

Iterative methods, and in particular Krylov subspace methods, are ubiquitous in scientific computing. Algorithms such as power iteration, Golub-Kahan Bidiagonalization, Arnoldi iteration, and the Lanczos iteration, are used in basic subroutines for matrix inversion, solving linear systems, linear programming, low-rank approximation, and numerous other fundamental linear algebra primitives [Saa81, LS13]. A common technique in the analysis of Krylov methods is the use of Chebyshev polynomials, which can be applied to the singular values of a matrix to implement an approximate interval or step function [MH02, Riv20]. Further, Chebyshev polynomials reduce the degree required to accurately approximate such functions, leading to significantly fewer iterations and faster running time. In this paper we investigate the power of Krylov methods for low-rank approximation in the matrix-vector product model.

The Matrix-Vector Product Model. In this model, there is an underlying matrix \mathbf{A} , which is often implicit, and for which the only access to \mathbf{A} is via matrix-vector products. Namely, the algorithm chooses a query vector v^1 , obtains the product $\mathbf{A} \cdot v^1$, chooses the next query vector v^2 , which is any randomized function of v^1 and $\mathbf{A} \cdot v^1$, then receives $\mathbf{A} \cdot v^2$, and so on. If \mathbf{A} is a non-symmetric matrix, we assume access to products of the form $\mathbf{A}^T v$ as well. We refer to the minimal number q of queries needed by the algorithm to solve a problem with constant probability as the *query complexity*. We note that upper bounds on the query complexity immediately translate to running time bounds for the RAM model, when \mathbf{A} is explicit, since a matrix-vector product can be implemented in $nnz(\mathbf{A})$ time, i.e., the number of non-zero entries in the matrix. Since this model captures a large family of iterative methods, it is natural to ask whether Krylov subspace based methods yield optimal algorithms, where the complexity measure of interest is the number of matrix-vector products.

This model and related vector-matrix-vector query models were formalized for a number of problems in [SWYZ19, RWZ20], though the model is standard for measuring efficiency in scientific computing and numerical linear algebra, see, e.g., [BFG96]; in that literature, methods that use only matrix-vector products are called *matrix-free*. Subsequently, for the problem of estimating the top eigenvector, nearly tight bounds were obtained in [SAR18, BHSW20]. Also, for the problem of estimating the trace of a positive semidefinite matrix, tight bounds were obtained in [MMMW21] (see, also [WWZ14], where tight bounds were shown in the related vector-matrix-vector query model). For recovering a planted clique from a random graph, upper and lower bounds were obtained in [RWYZ21]. In the non-adaptive setting, where v^1, \ldots, v^q , are chosen before making any queries to **A**, this is equivalent to the *sketching model*, which is thoroughly studied on its own (see, e.g., [Nel11, Woo14]), and in the context of data streams [Mut05, LNW14b].

Why is the matrix **A** implicit? A small query complexity q leads to an algorithm running in time $O(T(\mathbf{A}) \cdot q + P(n,d,q))$, where $T(\mathbf{A})$ is the time to multiply the $n \times d$ matrix **A** by an arbitrary vector, and P(n,d,q) is the time needed to form the queries and process the query responses, which is typically small. When the matrix **A** is given as a list of $nnz(\mathbf{A})$ non-zero entries, then $T(\mathbf{A}) \leq nnz(\mathbf{A})$. However, in many problems **A** is not given explicitly, and it is too expensive to write **A** down. Indeed, one may be given **A** but want to compute a low-rank approximation to the

"covariance" (Gram) matrix $\mathbf{A}^{\top}\mathbf{A}$, and computing $\mathbf{A}^{\top}\mathbf{A}$ is too slow [MW17a]. More generally, one may be given $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^{\top}$ and a function $f: \mathbb{R} \to \mathbb{R}$, and want to compute matrix-vector products with the generalized matrix function $f(\mathbf{A}) = \mathbf{U}f(\boldsymbol{\Sigma})\mathbf{V}^{\top}$, where \mathbf{U} has orthonormal columns, \mathbf{V}^{\top} has orthonormal rows, $\boldsymbol{\Sigma}$ is a diagonal matrix, and f is applied entry-wise to each entry on the diagonal.

The covariance matrix corresponds to $f(x) = x^2$, and other common functions f include the matrix exponential $f(x) = e^x$ and low-degree polynomials. For instance, when \mathbf{A} is the adjacency matrix of an undirected graph, $f(x) = x^3/6$ is used to count the number of triangles [Tso08, Avr10]. Yet another example is when \mathbf{A} is the Hessian \mathbf{H} of a neural network with a huge number of parameters, for which it is often impossible to compute or store the entire Hessian [GKX19]. Typically $\mathbf{H} \cdot v$, for any chosen vector v, is computed using Pearlmutter's trick [Pea94]. However, even with Pearlmutter's trick and distributed computation on modern GPUs, it takes 20 hours to compute the eigendensity of a single Hessian \mathbf{H} with respect to the cross-entropy loss on the CIFAR-10 dataset from a set of fixed weights for ResNet-18 [KH+09], which has approximately 11 million parameters [HZRS16, GKX19]. This time is directly proportional to the number of matrix-vector products, and therefore minimizing this quantity is crucial.

Algorithms and Lower Bounds for Low-Rank Approximation. The low-rank approximation problem is well studied in numerical linear algebra, with countless applications to clustering, data mining, principal component analysis, recommendation systems, and many more. (For surveys on low-rank approximation, see the monographs [KV09, Mah11, Woo14] and references therein.) In this problem, given an implicit $n \times d$ matrix \mathbf{A} , the goal is to output a matrix $\mathbf{Z} \in \mathbb{R}^{d \times k}$ with orthonormal columns such that

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{X} \leq (1 + \epsilon) \min_{\mathbf{U}: \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_{k}} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{X}, \tag{1.1}$$

where $\|\cdot\|_X$ denotes some norm. Note that given **Z**, one can compute **AZ** with an additional k queries, which will be negligible, and then $(\mathbf{AZ})\cdot\mathbf{Z}^{\top}$ is a rank-k matrix written in factored form, i.e., as the product of an $n\times k$ matrix and a $k\times d$ matrix. Among other things, low-rank approximation provides (1) a compression of **A** from nd parameters to (n+d)k parameters, (2) faster matrix-vector products, since $\mathbf{AZ}\cdot\mathbf{Z}^{\top}\cdot y$ can be computed in O((n+d)k) time for an arbitrary vector y, as opposed to the O(nd) time needed to compute $\mathbf{A}\cdot y$, and (3) de-noising, as often matrices **A** are close to low-rank (e.g., they are the product of latent factors) but only high rank due to noise.

Despite its tremendous importance, the optimal matrix-vector product complexity of low-rank approximation is unknown for any commonly used norm. The best known upper bound is due to Musco and Musco [MM15], who achieve $\tilde{O}(k/\epsilon^{1/2})$ queries for both the case when $\|\cdot\|_X$ is the commonly studied Frobenius norm $\|\mathbf{B}\|_F = \left(\sum_{i,j} \mathbf{B}_{i,j}^2\right)^{1/2}$ as well as when $\|\cdot\|_X$ is the Spectral (operator) norm $\|\mathbf{B}\|_2 = \sup_{\|y\|_2=1} \|\mathbf{B}y\|_2$.

On the lower bound front, there is a trivial lower bound of k, since **A** may be full rank and achieving (1.1) requires k matrix-vector products since one must reconstruct the column span of **A** exactly. However, no lower bounds in terms of the approximation factor ϵ were known. We note that

¹We let $\tilde{O}(f) = f \cdot \text{poly}(\log(dk/\epsilon))$.

Simchowitz, Alaoui and Recht [SAR18] prove lower bounds for approximating the top r eigenvalues of a symmetric matrix; however these guarantees are incomparable to those that follow from a low-rank approximation, even when the norm $\|\cdot\|_X$ is the operator norm (see Appendix A for a brief discussion).

Relationship to the Sketching Literature. Low-rank approximation has been extensively studied in the sketching literature which, when **A** is given explicitly, can achieve $O(\text{nnz}(\mathbf{A}))$ time both for the Frobenius norm [CW13, MM13a, NN13], as well as for Schatten-p norms [LW20]. However, these works require reading all of the entries in **A**, and thus do not apply to any of the settings mentioned above. Further, the matrix-vector query model is especially important for problems such as trace estimation, where a low-rank approximation is used to first reduce the variance [MMMW21]. As trace estimation is often applied to implicit matrices, e.g., in computing Stochastic Lanczos Quadrature (SLQ) for Hessian eigendensity estimation [GKX19], in studying the effects of batch normalization and residual connections in neural networks [YGKM20], and in computing a disentanglement regularizer for deep generative models [PPZ⁺20], sketching algorithms for low-rank approximation often do not apply.

Another important application is low-rank approximation of covariance matrices [MW17a], for which the covariance matrix is not given explicitly. Here, we have a data matrix \mathbf{A} and we want a low-rank approximation for $\mathbf{A}\mathbf{A}^{\top}$. Even when \mathbf{S} is a sparse sketching matrix, the matrix $\mathbf{S}\mathbf{A}$ is no longer sparse, and one needs to multiply $\mathbf{S}\mathbf{A}$ by \mathbf{A}^{\top} to obtain a sketch of $\mathbf{S}\mathbf{A}\mathbf{A}^{\top}$, which is a dense matrix-matrix multiplication. Moreover, when viewed in the matrix-vector product model, sketching algorithms obtain provably worse query complexity than existing iterative algorithms (see Table 1.1 for a comparison). Further, as modern GPUs often do not exploit sparsity, even when the matrix \mathbf{A} is given, a GPU may not be able to take advantage of sparse queries, which means the total time taken is proportional to the number of matrix-vector products.

Motivating Schatten-*p* Norms. The Schatten norms for $1 \le p < 2$ are more robust than the Frobenius norm, as they dampen the effect of large singular values. In particular, the Schatten-1 norm, also known as the nuclear norm, has been widely used for robust PCA [XCS10, CLMW11, YPCC16] as well as a convex relaxation of matrix rank in matrix completion [CR09, CP10], low-dimensional Euclidean embeddings [RFP10, TDSL00, RS00], image denoising [GZZF14, GXM+17] and tensor completion [YZ16]. In contrast, for p > 2, Schatten norms are more sensitive to large singular values and provide an approximation to the operator norm. In particular, for a rank r matrix, it is easy to see that setting $p = \log(r)/\eta$ yields a $(1 + \eta)$ -approximation to the operator norm (i.e., $p = \infty$). While the Block Krylov algorithm of Musco and Musco [MM15] implies a matrix-vector query upper bound of $\tilde{O}(k/\epsilon^{1/2})$ for Schatten-∞ low-rank approximation, the exact complexity of this problem remains an outstanding open problem. When p > 2, we can interpolate between Frobenius and operator norm, and setting p to be a large fixed constant can be a proxy for Schatten-∞ low-rank approximation, with significantly fewer matrix-vector products (see Theorem 5.1).

Our Central Question. The main question of our work is:

What is the matrix-vector product complexity of low-rank approximation for the Frobenius norm, and more

1.1 Our Results

Problem	Frobenius	Schatten- p , $p \in [1, 2)$	Schatten- p , $p > 2$
Sketching [CW09, LW20]	$\Theta(k/\epsilon)$	$\Omega(k^{2/p}/\epsilon^{4/p+1})$	$\Omega(\min(n,d)^{1-2/p})$
Block Krylov [MM15]	$\tilde{O}(k/\epsilon^{1/2})$	N/A	N/A
Our Upper Bound	$\tilde{O}(k/\epsilon^{1/3})$	$\tilde{O}(k/\epsilon^{1/3})$	$\tilde{O}(kp^{1/6}/\epsilon^{1/3})$
Our Lower Bound	$\Omega(1/\epsilon^{1/3})$	$\Omega(1/\epsilon^{1/3})$	$\Omega(1/\epsilon^{1/3})$

Figure 1.1: Prior Upper and Lower Bounds on the Matrix Vector Product Complexity for Frobenius and Schatten-p low-rank Approximation. The $\operatorname{poly}(k/\epsilon)$ factors in prior sketching work for Schatten-p are not explicit, but we have computed lower bounds on them to illustrate our improvements. Our bounds are optimal, up to logarithmic factors, for constant k. For $p > \log(d)/\epsilon$, spectral low-rank approximation [MM15] implies an $\tilde{O}(k/\sqrt{\epsilon})$ upper bound.

We begin by stating our results for Frobenius and more generally, Schatten-p norm low-rank approximation for any $p \ge 1$; see Table 1.1 for a summary.

Theorem 1.1 (Query Upper Bound, informal Theorem 5.1). Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, a target rank $k \in [d]$, an accuracy parameter $\epsilon \in (0,1)$ and any (not necessarily constant) $p \in [1, O(\log(d)/\epsilon)]$, there exists an algorithm that uses $\tilde{O}(kp^{1/6}/\epsilon^{1/3})$ matrix-vector products and outputs a $d \times k$ matrix \mathbf{Z} with orthonormal columns such that with probability at least 99/100,

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p} \leq (1 + \epsilon) \min_{\mathbf{U}: \ \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{\mathcal{S}_p}.$$

When $p \ge \log(d)/\epsilon$, we get $\tilde{O}(k/\sqrt{\epsilon})$ matrix-vector products.

We note that for Frobenius norm low-rank approximation (Schatten p for p=2), we improve the prior matrix-vector product bound of $\tilde{O}(k/\epsilon^{1/2})$ by Musco and Musco [MM15] to $\tilde{O}(k/\epsilon^{1/3})$. For Schatten-p low-rank approximation for $p\in[1,2)$, we improve work of Li and Woodruff [LW20] who require query complexity at least $\Omega(k^{2/p}/\epsilon^{4/p+1})$, which is a polynomial factor worse in both k and $1/\epsilon$ than our $\tilde{O}(k/\epsilon^{1/3})$ bound.

For p>2, [LW20] obtain a query complexity of $\Omega(\min(n,d)^{1-2/p})$. We drastically improve this to $\tilde{O}(k/\epsilon^{1/3})$, which does not depend on d or n at all. Setting $p=\log(d)/\epsilon$ suffices to obtain a $(1+\epsilon)$ -approximation to the spectral norm $(p=\infty)$, and we obtain an $\tilde{O}(k/\sqrt{\epsilon})$ query algorithm, matching the best known bounds for spectral low-rank approximation [MM15]. When $p>\log(d)/\epsilon$, we can simply run Block Krylov for $p=\infty$.

Remark 1.2 (Comments on the RAM Model). Although our focus is on minimizing the number of matrix-vector products, which is the key resource in the applications described above, our bounds also improve the running time of low-rank approximation algorithms when the matrix $\bf A$ has a small number of non-zero entries and is explicitly given. For simplicity, we state our bounds and those of previous work without using algorithms for fast matrix multiplication; similar improvements hold when using such algorithms. When $nnz(\bf A)=O(n)$, for Frobenius norm low-rank approximation, work in the sketching literature, and in particular [ACW17] (building off of [CW13, NN13, Coh16]), achieves $O(nk^2/\epsilon)$ time. In contrast, in this setting our runtime is $\tilde{O}(nk^2/\epsilon^{2/3})$. Similarly, for Schatten-p low-rank approximation for $p \in [1,2)$, the previous best [LW20] requires $\tilde{\Omega}(nk^{4/p}/\epsilon^{(8/p-2)})$ time, while for p > 2 [LW20] requires $\tilde{\Omega}(nd^{2(1-2/p)}(k/\epsilon)^{4/p})$ time. In both cases our runtime is only $\tilde{O}(nk^2p^{1/3}/\epsilon^{2/3})$. We obtain analogous improvements when the sparsity $nnz(\bf A)$ is allowed to be $n(k/\epsilon)^C$ for a small constant C > 0.

Next, we state our lower bounds on the matrix-vector query complexity of Schatten-p low-rank approximation.

Theorem 1.3 (Query Lower Bound for constant p, informal Theorem 6.1 and Theorem 6.4). Given $\varepsilon > 0$, and a fixed constant $p \ge 1$, there exists a distribution \mathcal{D} over $n \times n$ matrices such that for $\mathbf{A} \sim \mathcal{D}$, any algorithm that with at least constant probability outputs a unit vector v such that $\|\mathbf{A}(\mathbf{I} - vv^{\top})\|_{\mathcal{S}_p}^p \le (1 + \varepsilon) \min_{\|u\|_2 = 1} \|\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p}^p$ must perform $\Omega(1/\varepsilon^{1/3})$ matrix-vector queries to \mathbf{A} .

Remark 1.4. We note that this is the first lower bound as a function of ϵ for this problem, even for the well-studied case of p = 2, achieving an $\Omega(1/\epsilon^{1/3})$ bound, which is tight for any constant k, simultaneously for all constant $p \ge 1$.

Remark 1.5. Braverman, Hazan, Simchowitz and Woodworth [BHSW20] and Simchowitz, Alaoui and Recht [SAR18] establish eigenvalue estimation lower bounds that we use in our arguments, but their results do not directly imply low-rank approximation lower bounds for any matrix norm that we are aware of, including spectral low-rank approximation, i.e., $p = \infty$ (see Appendix A).

Matrix Polynomials and Streaming Algorithms. Since our algorithms are based on iterative methods, they generalize naturally to low-rank approximations of matrices of the form $\mathbf{A}(\mathbf{A}^{\mathsf{T}}\mathbf{A})^{\ell}$ and $(\mathbf{A}^{\mathsf{T}}\mathbf{A})^{\ell}$ for any integer ℓ , given \mathbf{A} as input. We defer the details to Appendix \mathbf{B} .

Since we work in the matrix-vector model, our algorithms naturally extend to the multi-pass turnstile streaming setting. Notably, for p>2, with $O(\log(d/\epsilon)p^{1/6}/\epsilon^{1/3})$ passes we are able to improve the $\tilde{O}\left(n\left(\frac{kn^{1-2/p}}{\epsilon^2}+\frac{k^{2/p}+n^{1-2/p}}{\epsilon^{2+2/p}}\right)\right)$ memory bound of [LW20] to $\tilde{O}\left(nk/\epsilon^{1/3}\right)$. We defer the details to Appendix C.

1.2 Open Questions

We note that our lower bounds are tight only when the target rank k and Schatten norm p are fixed constants. In particular, it is open to obtain matrix-vector lower bounds that grow as a function of k, p and $1/\epsilon$. For the important special case of Spectral low-rank approximation ($p = \infty$), it is open to obtain any lower bound that grows as a function of $1/\epsilon$, even when the target rank k = 1

(see Appendix A for more details). We also note that improving our upper bound to even $p^{1/6-o(1)}$ would imply a faster algorithm for Spectral low-rank approximation, addressing the main open question in [Woo14].

2 Technical Overview

For our technical overview, we drop polylogarithmic factors appearing in the analysis and assume the input **A** is a symmetric $n \times n$ matrix (we handle arbitrary $n \times d$ matrices in Section 5).

2.1 Algorithms for Low-Rank Approximation

We first describe our algorithm for the special case of rank-1 approximation in the Frobenius norm, i.e., p = 2. Our algorithm is inspired by the Block Krylov algorithm of Musco and Musco [MM15]. Briefly, their algorithm begins with a random starting vector g (block size is 1) and computes the Krylov subspace $\mathbf{K} = [\mathbf{A}g; \mathbf{A}^2g; \dots; \mathbf{A}^qg]$, for $q = O(1/\epsilon^{1/2})$. Next, their algorithm computes an orthonormal basis for the column span of \mathbf{K} , denoted by a matrix \mathbf{Q} , and outputs the top singular vector of $\mathbf{Q}^{\mathsf{T}}\mathbf{A}^2\mathbf{Q}$, denoted by z (see Algorithm 5.6 for a formal description). It follows from Theorem 1, guarantee (1) in [MM15] that

$$\left\| \mathbf{A} \left(\mathbf{I} - z z^{\mathsf{T}} \right) \right\|_{F}^{2} \le (1 + \epsilon) \min_{\|u\|_{2} = 1} \left\| \mathbf{A} \left(\mathbf{I} - u u^{\mathsf{T}} \right) \right\|_{F}^{2}, \tag{2.1}$$

and it is easy to see that this algorithm requires $\Theta\left(1/\epsilon^{1/2}\right)$ matrix-vector products. A naïve analysis requires an $O(1/\epsilon)$ -degree polynomial in the matrix **A** to obtain (2.1), while [MM15] use Chebyshev polynomials to approximate the threshold function between first and second singular value, and save a quadratic factor in the degree. The guarantee in (2.1) then follows from observing that the best vector in the Krylov subspace is at least as good as the one that exists using Chebyshev polynomial approximation.

Algorithm 2.1 (Algorithm Sketch for Frobenius rank-1 LRA).

Input: An $n \times n$ symmetric matrix **A**, accuracy parameter $0 < \varepsilon < 1$.

- 1. Run Block Krylov for $O(1/\epsilon^{1/3})$ iterations with a random starting vector g. Let z_1 be the resulting output.
- 2. Run Block Krylov for $O(\log(n/\epsilon))$ iterations, but initialize with an $n \times b$ random matrix **G**, where $b = O(1/\epsilon^{1/3})$. Let z_2 be the resulting output.

Output:
$$z = \arg \max_{z_1, z_2} (\|\mathbf{A}z_1\|_2^2, \|\mathbf{A}z_2\|_2^2).$$

Our starting point is the observation that while we require degree $\Theta(1/\epsilon^{1/2})$ to separate the first and second singular values, if any subsequent singular value is sufficiently separated from σ_1 , a significantly smaller degree polynomial suffices. In the context of Krylov methods, this translates

to the intuition that starting with a matrix G with b columns (block size is b) should result in fewer iterations to find some vector in the top b subspace of G. On the other hand, if no such singular value exists, the norm of the tail must be large and we can get away with a less accurate solution. We show that we can indeed exploit this trade-off by running Block Krylov on two different scales in parallel and then combine the solution. In particular, we use Algorithm 2.1.

Algorithm 2.1 captures the extreme points of the trade-off between the size of the starting matrix and the number of iterations, such that the total number of matrix-vector products is at most $\tilde{O}(1/\epsilon^{1/3})$. Further, we can compute the squared Euclidean norms of Az_1 and Az_2 with an additional matrix-vector product, and it remains to analyze the Frobenius cost of projecting A on the subspace $I - zz^T$, where z is the unit vector output by Algorithm 2.1.

Using gap-independent guarantees for Block Krylov (see Lemma 5.2 for a formal statement), it follows that with $O(1/\epsilon^{1/3})$ iterations, we have

$$\|\mathbf{A}z_1\|_2^2 \ge \sigma_1^2(\mathbf{A}) - \epsilon^{2/3}\sigma_2^2(\mathbf{A}).$$
 (2.2)

In contrast, using gap-dependent guarantees (see Lemma 5.4) for Block Krylov initialized with block size b, it follows that for any $\gamma > 0$, running $q = \log(1/\gamma) \cdot \sqrt{\sigma_1(\mathbf{A})/(\sigma_1(\mathbf{A}) - \sigma_b(\mathbf{A}))}$ iterations results in z_2 such that

$$\|\mathbf{A}z_2\|_2^2 \ge \sigma_1^2(\mathbf{A}) - \gamma \sigma_2^2(\mathbf{A}).$$
 (2.3)

If $\sigma_b(\mathbf{A}) \leq \sigma_1(\mathbf{A})/2$, we can set $\gamma = \epsilon/n$ in Equation (2.3) to obtain a highly accurate solution. Further, regardless of the input instance, Step 3 in Algorithm 2.1 ensures that we get the best of both guarantees, (2.2) and (2.3). Then, observing that $\mathbf{I} - zz^{\mathsf{T}}$ is an orthogonal projection matrix (see Definition 4.1) and using the Pythagorean Theorem for Euclidean space we have:

$$\|\mathbf{A} (\mathbf{I} - zz^{\mathsf{T}})\|_F^2 = \|\mathbf{A}\|_F^2 - \|\mathbf{A}zz^{\mathsf{T}}\|_F^2 = \|\mathbf{A}\|_F^2 - \|\mathbf{A}z\|_2^2,$$
 (2.4)

where the second inequality follows from unitary invariance (see Fact 4.8) of the Frobenius norm and that the squared Frobenius norm of a rank-1 matrix $\mathbf{A}z$ (vector) is equal to its squared Euclidean norm. If it happens that $\sigma_2(\mathbf{A}) \leq \sigma_1(\mathbf{A})/2$, i.e., a constant gap exists between the first two singular values, then since guarantee (2.3) implies that $\|\mathbf{A}z\|_2^2 \geq \sigma_1^2(\mathbf{A}) - (\epsilon/n)\sigma_2^2(\mathbf{A})$, we can plug this into (2.4) to yield a $(1+\epsilon/n)$ -approximate solution. Hence, we focus on instances where $\sigma_2(\mathbf{A}) > \sigma_1(\mathbf{A})/2$.

Consider the case where the Frobenius norm of the tail is large, i.e., $\|\mathbf{A} - \mathbf{A}_1\|_F^2 \ge \sigma_2^2(\mathbf{A})/\epsilon^{1/3}$, where \mathbf{A}_1 is the best rank-1 approximation to \mathbf{A} . Then we only require an $\epsilon^{2/3}$ -approximate solution (plugging guarantee (2.2) into (2.4)) since

$$\|\mathbf{A} \left(\mathbf{I} - z_1 z_1^{\mathsf{T}}\right)\|_F^2 \le \|\mathbf{A}\|_F^2 - \sigma_1^2(\mathbf{A}) + \epsilon^{2/3} \sigma_2^2(\mathbf{A}) \le \|\mathbf{A} - \mathbf{A}_1\|_F^2 + \epsilon \|\mathbf{A} - \mathbf{A}_1\|_F^2.$$
 (2.5)

Otherwise, $\sum_{i=2}^n \sigma_i^2(\mathbf{A}) < \sigma_2^2(\mathbf{A})/\epsilon^{1/3}$, which implies that there is a constant gap between the second and b-th singular values, where $b = O(1/\epsilon^{1/3})$. To see this, observe if $\sigma_b(\mathbf{A}) > \sigma_2(\mathbf{A})/4$, then $\sum_{i=2}^n \sigma_i^2(\mathbf{A}) \ge \sum_{i=2}^b \sigma_i^2(\mathbf{A}) \ge b\sigma_2^2(\mathbf{A})/4$, which is a contradiction when $b > 10/\epsilon^{1/3}$, and thus $\sigma_b(\mathbf{A}) \le \sigma_2(\mathbf{A})/4 < \sigma_1/2$. Now we can apply guarantee (2.3) with $q = O(\log(n/\epsilon))$ and conclude $\|\mathbf{A}z\|_2^2 \ge \sigma_1^2(\mathbf{A}) - (\epsilon/n)\sigma_2^2(\mathbf{A})$, yielding a highly accurate solution yet again. Overall, this suffices to obtain a $(1 + \epsilon)$ -approximate solution with $\tilde{O}(1/\epsilon^{1/3})$ matrix-vector queries.

Challenges in generalizing to Schatten $p \neq 2$ and rank k > 1. The outline above crucially relies on the norm of interest being Frobenius. In particular, we use the Pythagorean Theorem to analyze the cost of the candidate solution in Equation (2.4); however, the Pythagorean Theorem does not hold for non-Euclidean spaces. Therefore, a priori, it is unclear how to analyze the Schatten-p norm of a candidate rank-1 approximation. A proxy for the Pythagorean Theorem that holds for Schatten-p norms is Mahler's operator inequality (see Fact 4.11), which is in the right direction but holds only for $p \geq 2$, whereas we would like to handle all $p \geq 1$. Separately, for p > 2, the case where the tail is small corresponds to $\|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p \leq \sigma_2^p(\mathbf{A})/\epsilon^{1/3}$. Therefore, naïvely extending the above argument requires picking a block size that scales proportional to $O(2^p/\epsilon^{1/3})$ to induce a constant gap between σ_1 and σ_b , and the number of matrix-vector products scales exponentially in p.

Finally, in the above outline, we also crucially use that $\|\mathbf{A}zz^{\top}\|_{F}^{2} = \|\mathbf{A}z\|_{2}^{2}$. Observe that this no longer holds if we replace z with a matrix \mathbf{Z} that has k orthonormal columns. Therefore, it remains unclear how to relate $\|\mathbf{A}\mathbf{Z}\|_{\mathcal{S}_{p}}^{p}$ to $\|\mathbf{A}\mathbf{Z}_{*,i}\|_{2}^{2}$, yet the vector-by-vector error guarantee obtained by Block Krylov (see Lemmas 5.2 and 5.4) only bounds the latter.

Handling all Schatten-p **Norms and** k > 1. We modify our algorithm to run Block Krylov on \mathbf{A}^{\top} and obtain an orthonormal matrix \mathbf{W} such that for all $i \in [k]$,

$$\left\|\mathbf{A}^{\mathsf{T}}\mathbf{W}_{*,i}\right\|^{2} \ge \sigma_{i}^{2}(\mathbf{A}) - \gamma \sigma_{k+1}^{2}(\mathbf{A}),\tag{2.6}$$

for some $\gamma > 0$. We then analyze the cost $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p}^p$, where \mathbf{Z} is a basis for $\mathbf{A}^{\top}\mathbf{W}$. Our key insight is to interpret the input matrix \mathbf{A} as a partitioned operator (block matrix) and invoke *pinching inequalities* for such operators. Pinching inequalities were originally introduced to understand unitarily invariant norms over direct sums of Hilbert spaces [VN37, Sch60]. In our setting, given a block matrix $\mathbf{M} = \begin{pmatrix} \mathbf{M}^{(1)} & \mathbf{M}^{(2)} \\ \mathbf{M}^{(3)} & \mathbf{M}^{(4)} \end{pmatrix}$, the *pinching inequality* (see Fact 4.13) implies that for all $p \geq 1$,

$$\|\mathbf{M}\|_{\mathcal{S}_p}^p \ge \|\mathbf{M}^{(1)}\|_{\mathcal{S}_p}^p + \|\mathbf{M}^{(4)}\|_{\mathcal{S}_p}^p.$$
 (2.7)

A priori, it is unclear how to use Equation (2.7) to bound $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p}^p$. First, we establish a general inequality for the Schatten norm of a matrix times an orthogonal projection. Let **P** and **Q** be any $n \times n$ orthogonal projection matrices with rank k (see Definition 4.1). Then, we prove (see Lemma 5.5 for details) that for any matrix **A**,

$$\|\mathbf{A}\|_{\mathcal{S}_p}^p \ge \|\mathbf{P}\mathbf{A}\mathbf{Q}\|_{\mathcal{S}_p}^p + \|(\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_p}^p.$$
(2.8)

To obtain this inequality, we use a rotation argument along with the fact that the Schatten-p norms are unitarily invariant to show that $\|\mathbf{A}\|_{\mathcal{S}_p}^p = \left\| \begin{pmatrix} \mathbf{A}^{(1)} & \mathbf{A}^{(2)} \\ \mathbf{A}^{(3)} & \mathbf{A}^{(4)} \end{pmatrix} \right\|_{\mathcal{S}_p}^p$, where $\|\mathbf{A}^{(1)}\|_{\mathcal{S}_p} = \|\mathbf{P}\mathbf{A}\mathbf{Q}\|_{\mathcal{S}_p}$ and $\|\mathbf{A}^{(4)}\|_{\mathcal{S}_n} = \|(\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_p}$, and then we can apply Equation (2.7) to the block matrix above.

Once we have established Equation (2.8), we can set $P = WW^{\top}$ and set $Q = ZZ^{\top}$ to be the projection matrix corresponding to the column span of $A^{\top}WW^{\top}$. Then, we have that $PAQ = WW^{\top}A$ and $(I - P) A (I - Q) = A (I - ZZ^{\top})$, and combined with (2.8) this yields

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\mathsf{T}} \right) \right\|_{\mathcal{S}_{p}}^{p} \leq \left\| \mathbf{A} \right\|_{\mathcal{S}_{p}}^{p} - \left\| \mathbf{W} \mathbf{W}^{\mathsf{T}} \mathbf{A} \right\|_{\mathcal{S}_{p}}^{p}. \tag{2.9}$$

To obtain a bound on $\|\mathbf{W}\mathbf{W}^{\mathsf{T}}\mathbf{A}\|_{\mathcal{S}_p}^p$, we appeal to the per-vector guarantees in Equation (2.6). However, translating from ℓ_2^2 error to $\sigma_p^p(\mathbf{W}^{\mathsf{T}}\mathbf{A})$ incurs a mixed guarantee (see Lemma 5.7 for details):

$$\left\|\mathbf{W}\mathbf{W}^{\top}\mathbf{A}\right\|_{\mathcal{S}_{p}}^{p} \geq \left\|\mathbf{A}_{k}\right\|_{\mathcal{S}_{p}}^{p} - O(\gamma p) \sum_{i \in [k]} \sigma_{k+1}^{2}\left(\mathbf{A}\right) \sigma_{i}^{p-2}\left(\mathbf{A}\right).$$

To use this bound, we require $\sigma_1(\mathbf{A})$ to be comparable to $\sigma_{k+1}(\mathbf{A})$ and thus we require an involved case analysis, which appears in the proof of Theorem 5.1.

Avoiding an exponential dependence on p. Our main insight here is that we do not require a block size that induces a constant gap between singular values. Instead, we first observe that if the block size b is large enough such that $\sigma_b \leq \sigma_2/(1+1/p)$, then $O(\log(n/\epsilon)\sqrt{p})$ iterations suffice to obtain a vector z such that $\|\mathbf{A}z\|_2^2 \geq \sigma_1^2(\mathbf{A}) - (\epsilon/n)\sigma_2^2(\mathbf{A})$. Therefore, we can trade-off the threshold for the Schatten norm of the tail with the number of iterations as follows: if $\|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p \leq \frac{1}{p^{1/3}\epsilon^{1/3}}\sigma_2^p(\mathbf{A})$, then setting $b = (1+1/p)^p/(\epsilon p)^{1/3} = \Theta(1/(\epsilon p)^{1/3})$ suffices to induce a gap of 1+1/p with block size b. The total number of matrix-vector products is $O(b \cdot \log(n/\epsilon)\sqrt{p}) = \tilde{O}(p^{1/6}/\epsilon^{1/3})$, since p can be assumed to be at most $(\log n)/\epsilon$. Otherwise, $\|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p > \frac{1}{p^{1/3}\epsilon^{1/3}}\sigma_2^p(\mathbf{A})$, and we only require a $(1+\epsilon^{2/3}/p^{1/3})$ -approximate solution instead (compare with Equation (2.5)). Using gap-independent bounds (see Lemma 5.2), it suffices to start with block size 1 and run $O(\log(n/\epsilon)p^{1/6}/\epsilon^{1/3})$ iterations to obtain a $(1+\epsilon^{2/3}/p^{1/3})$ -approximate solution.

Avoiding a Gap-Dependent Bound. We note that even when there is a constant gap between the first and second singular values, and the per vector guarantee is highly accurate, i.e., for all $i \in [k]$, $\|\mathbf{AZ}_{*,i}\|^2 \ge \sigma_i^2(\mathbf{A}) - \operatorname{poly}\left(\frac{\epsilon}{d}\right)\sigma_{k+1}^2(\mathbf{A})$, it is not clear how to lower bound $\|\mathbf{AZ}\|_{\mathcal{S}_p}^p$ in Equation 2.9. In general, the best bound we can obtain using the above equation is

$$\|\mathbf{AZ}\|_{\mathcal{S}_p}^p \ge \|\mathbf{A}_k\|_{\mathcal{S}_p}^p - O\left(\frac{\epsilon}{\text{poly}(d)}\right) \sigma_{k+1}^2 \cdot \sum_{i \in [k]} \sigma_i^{p-2}, \tag{2.10}$$

which may be vacuous when the top k singular values are significantly larger than σ_{k+1} and p > 2. One could revert to a gap-dependent bound, where the error is in terms of the gap between σ_1 and σ_{k+1} , which one could account for by running an extra factor of $O(\log(\sigma_1/\sigma_{k+1}))$ iterations.

To avoid this gap-dependent bound, we split **A** into a head part \mathbf{A}_H and a tail part \mathbf{A}_T , such that \mathbf{A}_H has all singular values that are at least $(1+1/d) \sigma_{k+1}$ and \mathbf{A}_T has the remaining singular values. We then bound $\|\mathbf{A}_H (\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)\|_{\mathcal{S}_p}$ and $\|\mathbf{A}_T (\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)\|_{\mathcal{S}_p}$ separately. Repeating the above analysis, we can obtain Equation (2.10) for \mathbf{A}_T instead, and since all singular values larger than σ_{k+1} in \mathbf{A}_T

are bounded, we can obtain $\|\mathbf{A}_T(\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)\|_{\mathcal{S}_p}^p \leq O(\epsilon k/\text{poly}(d)) \sigma_{k+1}^p$. To adapt the analysis for \mathbf{A}_T and obtain this bound, we use Cauchy's interlacing theorem to relate the j-th singular value of $\mathbf{A}_T(\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)$ to the $(i^* + j)$ -th singular value of $\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)$, where i^* is the rank of \mathbf{A}_H . We lower bound the $(i^* + j)$ -th singular value of $\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^\top)$ using the per vector guarantee of [MM15].

To bound $\|\mathbf{A}_H(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\mathsf{T}})\|_{S_n}$, we observe it has rank at most k and thus

$$\left\|\mathbf{A}_{H}\left(\mathbf{I}-\mathbf{Z}\mathbf{Z}^{\top}\right)\right\|_{\mathcal{S}_{p}} \leq \sqrt{k} \cdot \left\|\mathbf{A}_{H}\left(\mathbf{I}-\mathbf{Z}\mathbf{Z}^{\top}\right)\right\|_{F} = \sqrt{k} \cdot \sqrt{\left\|\mathbf{A}_{H}\right\|_{F}^{2} - \left\|\mathbf{A}_{H}\mathbf{Z}\right\|_{F}^{2}},$$

and we show how to bound this term in Section 5. Intuitively, while the k-dimensional subspace that we find can "swap out" singular vectors corresponding to singular values σ_i for which σ_i is very close to σ_{k+1} , since they serve equally well for a Schatten-p low-rank approximation, for singular values σ_i that are a bit larger than σ_{k+1} , the k-dimensional subspace we find cannot do this. More precisely, if y is a singular vector of \mathbf{A}_H with singular value σ_i , then the projection of y onto the k-dimensional subspace that our algorithm finds (namely, \mathbf{Z}) must be at least $1 - \sigma_{k+1}^2/((\sigma_i^2 - \sigma_{k+1}^2)\operatorname{poly}(d))$, which suffices to bound the above since the additive error is inversely proportional to σ_i^2 when $\sigma_i^2 \gg \sigma_{k+1}^2$, and so the very tiny additive error negates the effect of very large singular values.

2.2 Query Lower Bounds.

Our lower bounds rely on the hardness of estimating the smallest eigenvalue of a Wishart ensemble (see Definition 4.15), as established in recent work of Braverman, Hazan, Simchowitz and Woodworth [BHSW20]. In particular, [BHSW20] show that for a $d \times d$ instance \mathbf{W} of a Wishart ensemble, estimating $\lambda_d(\mathbf{W})$ (minimum eigenvalue) to additive error $1/d^2$ requires $\Omega(d)$ adaptive matrix-vector product queries (see Theorem 3.1 in [BHSW20]). To obtain hardness for Schatten-p low-rank approximation, we show that when $d = \Theta\left(1/\epsilon^{1/3}\right)$, any candidate unit vector z that satisfies $\|(\mathbf{I} - \mathbf{W}/5) (\mathbf{I} - zz^{\top})\|_{\mathcal{S}_p}^p \leq (1+\epsilon) \min_{\|u\|_2=1} \|(\mathbf{I} - \mathbf{W}/5) (\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p}^p$, can be used to obtain an estimate $\hat{\lambda}_d = \frac{5}{p} \left(1 - \|(\mathbf{I} - \mathbf{W}/5)z\|_2^p\right)$ such that $\hat{\lambda}_d = (1\pm 1/d^2)\lambda_d (\mathbf{I} - \mathbf{W}/5)$. Let $\mathbf{A} = (\mathbf{I} - \mathbf{W}/5)$. To show our query lower bound, in contrast to the analysis of our algorithm, the challenge is now to lower bound $\|\mathbf{A}(\mathbf{I} - zz^{\top})\|_{\mathcal{S}_p}^p$ in terms of $\|\mathbf{A}\|_{\mathcal{S}_p}^p$ and $\|\mathbf{A}z\|_2^p$ (contrast with Equation (2.9)).

Projection Cost via Araki-Lieb-Thirring. First, we note that the case of p = 2 is easy given the Pythagorean theorem. For $p \in [1,2)$, we can establish an inequality fairly straightforwardly: using the trace inner product definition of Schatten-p (see Definition 4.7) norms, we have,

$$\|\mathbf{A} \left(\mathbf{I} - zz^{\mathsf{T}}\right)\|_{\mathcal{S}_p}^p = \operatorname{Tr}\left(\left(\left(\mathbf{I} - zz^{\mathsf{T}}\right)^2 \mathbf{A}^2 \left(\mathbf{I} - zz^{\mathsf{T}}\right)^2\right)^{p/2}\right),\tag{2.11}$$

Since $p/2 \in [1/2, 1)$, we can use the reverse *Araki-Lieb-Thirring* inequality (see Fact 4.10) to show that

$$\operatorname{Tr}\left(\left(\left(\mathbf{I}-zz^{\top}\right)^{2}\mathbf{A}^{2}\left(\mathbf{I}-zz^{\top}\right)^{2}\right)^{p/2}\right) \geq \operatorname{Tr}\left(\left(\mathbf{I}-zz^{\top}\right)\mathbf{A}^{p}\left(\mathbf{I}-zz^{\top}\right)\right)$$

$$= \operatorname{Tr}\left(\mathbf{A}^{p}\right) - \operatorname{Tr}\left(\left(zz^{\top}\right)^{p/2}\left(\mathbf{A}^{2}\right)^{p/2}\left(zz^{\top}\right)^{p/2}\right)$$

$$\geq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}zz^{\top}\|_{\mathcal{S}_{p}}^{p}$$

$$(2.12)$$

where we use the cyclicity of the trace and again use reverse *Araki-Lieb-Thirring* (Fact 4.10) to show that

 $\operatorname{Tr}\left(\left(zz^{\top}\right)^{\frac{p}{2}}\left(\mathbf{A}^{2}\right)^{\frac{p}{2}}\left(zz^{\top}\right)^{\frac{p}{2}}\right) \leq \operatorname{Tr}\left(\left(zz^{\top}\mathbf{A}^{2}zz^{\top}\right)^{p/2}\right) = \left\|\mathbf{A}zz^{\top}\right\|_{\mathcal{S}_{v}}^{p}.$

Since we have $\|\mathbf{A}zz^{\top}\|_{\mathcal{S}_p}^p = \|\mathbf{A}z\|_2^p$, we conclude $\|\mathbf{A}(\mathbf{I} - zz^{\top})\|_{\mathcal{S}_p}^p \ge \|\mathbf{A}\|_{\mathcal{S}_p}^p - \|\mathbf{A}zz^{\top}\|_2^p$. This approach only works for $p \in [1,2)$; for p > 2 the application of *Araki-Lieb-Thirring* is reversed in Equation 2.12 (since p/2 > 1, see Fact 4.10) and we no longer get a lower bound on the cost in Equation 2.11. We therefore require a new approach.

Projection Cost via Norm Compression. Recall, z is the unit vector output by our candidate low-rank approximation and let $y = \mathbf{A}z/\|\mathbf{A}z\|_2$. We yet again interpret the input matrix \mathbf{A} as a partitioned operator by considering the projection of \mathbf{A} onto zz^{T} , yy^{T} and the projection away from these rank-1 subspaces. In particular, let $\mathbf{I} - yy^{\mathsf{T}} = \mathbf{Y}\mathbf{Y}^{\mathsf{T}}$, and $\mathbf{I} - zz^{\mathsf{T}} = \mathbf{Z}\mathbf{Z}^{\mathsf{T}}$, where \mathbf{Y} and \mathbf{Z} have orthonormal columns. Then, using a rotation argument, we show that

$$\|\mathbf{A}\|_{\mathcal{S}_p} = \left\| \begin{pmatrix} y^{\mathsf{T}} \mathbf{A} z & y^{\mathsf{T}} \mathbf{A} \mathbf{Z} \\ \mathbf{Y}^{\mathsf{T}} \mathbf{A} z & \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \end{pmatrix} \right\|_{\mathcal{S}_p}.$$

We define the *p*-compression of **A**, $C_{\mathbf{A},p}$:

$$\mathbf{C}_{\mathbf{A},p} = \begin{pmatrix} \| \mathbf{y}^{\mathsf{T}} \mathbf{A} \mathbf{z} \|_{\mathcal{S}_p} & \| \mathbf{y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \|_{\mathcal{S}_p} \\ \| \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{z} \|_{\mathcal{S}_p} & \| \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \|_{\mathcal{S}_p} \end{pmatrix}.$$

To relate the norms of **A** and $C_{A,p}$, we consider Audenaert's Norm Compression Conjecture [Aud08], a question in functional analysis concerning operator inequalities (see also [AK12]):

Conjecture 2.2 (Schatten-p Norm Compression). Let \mathbf{M} be a partitioned operator (block matrix) such that $\mathbf{M} = \begin{pmatrix} \mathbf{M}_1 & \mathbf{M}_2 \\ \mathbf{M}_3 & \mathbf{M}_4 \end{pmatrix}$. Let $\mathbf{C}_{\mathbf{M},p} = \begin{pmatrix} \|\mathbf{M}_1\|_{\mathcal{S}_p} & \|\mathbf{M}_2\|_{\mathcal{S}_p} \\ \|\mathbf{M}_3\|_{\mathcal{S}_p} & \|\mathbf{M}_4\|_{\mathcal{S}_p} \end{pmatrix}$ be a 2×2 matrix that denotes the Schatten-p compression of \mathbf{M} for any $p \geq 1$. Then, $\|\mathbf{M}\|_{\mathcal{S}_p} \geq \|\mathbf{C}_{\mathbf{M},p}\|_{\mathcal{S}_p}$ if $1 \leq p \leq 2$, and $\|\mathbf{M}\|_{\mathcal{S}_p} \leq \|\mathbf{C}_{\mathbf{M},p}\|_{\mathcal{S}_p}$ if $2 \leq p < \infty$.

We could simply appeal to this conjecture to obtain that for all p > 2,

$$\|\mathbf{A}\|_{\mathcal{S}_{p}} \leq \|\mathbf{C}_{\mathbf{A},p}\|_{\mathcal{S}_{p}} = \left\| \begin{pmatrix} \|yy^{\mathsf{T}}\mathbf{A}zz^{\mathsf{T}}\|_{\mathcal{S}_{p}} & \|yy^{\mathsf{T}}\mathbf{A}(I-zz^{\mathsf{T}})\|_{\mathcal{S}_{p}} \\ \|(\mathbf{I}-yy^{\mathsf{T}})\mathbf{A}zz^{\mathsf{T}}\|_{\mathcal{S}_{p}} & \|(\mathbf{I}-yy^{\mathsf{T}})\mathbf{A}(\mathbf{I}-zz^{\mathsf{T}})\|_{\mathcal{S}_{p}} \end{pmatrix} \right\|_{\mathcal{S}_{p}}.$$
 (2.13)

However, for our choice of y, $\|yy^{\top}\mathbf{A}(I-zz^{\top})\|_{S_p} = 0$. With padding and rotation arguments, we can then reduce our problem to a block matrix where the blocks in each row are aligned, i.e., each row is a scalar multiple of a fixed matrix (see Lemma 6.6). Then, we can use one of the few special cases of Conjecture 2.2 for aligned operators which has actually been proved, and appears in Fact 4.14. We can thus unconditionally obtain the inequality in Equation (2.13).

Now that we have reduced to the case where we have a 2×2 matrix with 3 non-zero entries, we would like to bound its Schatten-p norm. We explicitly compute the singular values of $C_{A,p}$ (see Fact 6.7), and then use the structure of the instance to directly lower bound $\|Az\|_2^p$ as follows:

$$\|\mathbf{A}z\|_{2}^{p} + \left(1 + O(\epsilon^{2p/3})\right)\|\mathbf{A} - \mathbf{A}_{1}\|_{\mathcal{S}_{p}}^{p} \ge \|\mathbf{C}_{\mathbf{A},p}\|_{\mathcal{S}_{p}}^{p} \ge \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p},$$
 (2.14)

where the last inequality follows from Equation (2.13). Since we understand the spectrum of the matrix **A**, we can explicitly compute all the terms in (2.14) above and show that we can obtain an accurate estimate of the minimum singular value of **A** from $\|\mathbf{A}z\|_2^p$. See details in Section 6.2.

3 Additional Related Work

Existing approaches to solve low-rank approximation problems under several norms fall into two broad categories: iterative methods and linear sketching. Iterative methods, such as Krylov subspace based methods, are captured by the matrix-vector product framework, whereas linear sketching allows for the choice of a matrix $\mathbf{S} \in \mathbb{R}^{t \times n}$, where t is the number of "queries", and then observes the product $\mathbf{S} \cdot \mathbf{A}$ and so on (see [Woo14] and references therein). The model has important applications to streaming and distributed algorithms and several recent works have focused on estimating spectral norms and the top singular values [AN13, LNW14a, LW16b, BBK+21], estimating Schatten and Ky-Fan norms [LW16b, LW17, LW16a, BKKS19] and low-rank approximation [CW13, MM13b, NN13, BDN15, Coh16].

In addition to studying unitarily invariant norms, such as the Schatten norm, there also has been significant amount of work on studying low-rank approximation under matrix ℓ_p norms [SWZ17, BBB+19, SWZ20, MW21] and weighted low-rank approximation [SJ03, RSW16, BWZ19], settings in which the problem is known to be NP-Hard. Finally, there has been a recent flurry of work on sublinear time algorithms for low-rank approximation under various structural assumptions on the input [MW17b, BW18, IVWW19, SW19, BCW20] and in quantum-inspired models [KP16, CLW18, Tan19, RSML18, GLT18, GSLW19, CCHW20].

4 Preliminaries

Given an $n \times d$ matrix **A** with rank r, and $n \ge d$, we can compute its singular value decomposition, denoted by $SVD(\mathbf{A}) = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\mathsf{T}}$, such that **U** is an $n \times r$ matrix with orthonormal columns, \mathbf{V}^{T} is an $r \times d$ matrix with orthonormal rows and $\mathbf{\Sigma}$ is an $r \times r$ diagonal matrix. The entries along the diagonal are the singular values of **A**, denoted by $\sigma_1, \sigma_2 \dots \sigma_r$. Given an integer $k \le r$, we define the truncated singular value decomposition of **A** that zeros out all but the top k singular values of **A**, i.e.,

 $\mathbf{A}_k = \mathbf{U} \mathbf{\Sigma}_k \mathbf{V}^{\mathsf{T}}$, where $\mathbf{\Sigma}_k$ has only k non-zero entries along the diagonal. It is well-known that the truncated SVD computes the best rank-k approximation to \mathbf{A} under any unitarily invariant norm, but in particular for any Schatten-p norm (defined below), we have $\mathbf{A}_k = \min_{\mathrm{rank}(\mathbf{X})=k} \|\mathbf{A} - \mathbf{X}\|_{\mathcal{S}_p}$. More generally, for any matrix \mathbf{M} , we use the notation \mathbf{M}_k and $\mathbf{M}_{\setminus k}$ to denote the first k components and all but the first k components respectively. We use $\mathbf{M}_{i,*}$ and $\mathbf{M}_{*,j}$ to refer to the i^{th} row and j^{th} column of \mathbf{M} respectively.

We use the notation I_k to denote a *truncated identity matrix*, that is, a square matrix with its top k diagonal entries equal to one, and all other entries zero. The dimension of I_k will be determined by context.

Definition 4.1 (Orthogonal Projection Matrices). Given a $d \times d$ symmetric matrix **P** and $k \in [d]$, **P** is a rank-k orthogonal projection matrix if rank(**P**) = k and **P**² = **P**.

It follows from the above definition that **P** has eigenvalues that are either 0 or 1 and admits a singular value decomposition of the form $\mathbf{U}\mathbf{U}^{\mathsf{T}}$ where **U** has k orthonormal columns.

Definition 4.2 (Unitary Matrices). Given a symmetric matrix $\mathbf{U} \in \mathbb{R}^{d \times d}$ we say \mathbf{U} is a unitary matrix if $\mathbf{U}^{\mathsf{T}}\mathbf{U} = \mathbf{U}\mathbf{U}^{\mathsf{T}} = \mathbf{I}$.

Definition 4.3 (Rotation Matrices). Given a symmetric matrix $\mathbf{R} \in \mathbb{R}^{d \times d}$ we say \mathbf{R} is a rotation matrix if \mathbf{R} is unitary and $\det(\mathbf{R}) = 1$.

Fact 4.4 (Courant-Fischer for Singular Values). *Given an* $n \times d$ *matrix* **A** *with singular values* $\sigma_1 \ge \sigma_2 \ge \ldots \ge \sigma_d$, *the following holds: for all* $i \in [d]$,

$$\sigma_i = \max_{S: \ dim(S)=i} \ \min_{x \in S: \ \|x\|_2=1} \ \left\|x^\top \mathbf{A}\right\|_2.$$

Fact 4.5 (Weyl's Inequality for Singular Values (see Exercise 22 [Tao20])). *Given* $n \times d$ *matrices* X, Y, *for any* i, $(j-1) \in [d]$ *such that* $i+j \leq d$,

$$\sigma_{i+i}(\mathbf{X} + \mathbf{Y}) \leq \sigma_i(\mathbf{X}) + \sigma_{i+1}(\mathbf{Y}).$$

Fact 4.6 (Bernoulli's Inequality). For any $x, p \in \mathbb{R}$ such that $x \ge -1$ and $p \ge 1$, $(1 + x)^p \ge 1 + px$.

Schatten Norms and Trace Inequalities. We recall some basic facts for Schatten-*p* norms. We also require the following trace and operator inequalities.

Definition 4.7 (Schatten-p Norm). Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, let $\sigma_1 \geq \sigma_2 \geq \ldots \geq \sigma_d$ be the singular values of \mathbf{A} . Then, for any $p \in [0, \infty)$, the Schatten-p norm of \mathbf{A} is defined as

$$\|\mathbf{A}\|_{\mathcal{S}_p} = \operatorname{Tr}\left(\left(\mathbf{A}^{\top}\mathbf{A}\right)^{p/2}\right)^{1/p} = \left(\sum_{i \in [d]} \sigma_i^p(\mathbf{A})\right)^{1/p}.$$

Fact 4.8 (Schatten-p norms are Unitarily Invariant). Given an $n \times d$ matrix \mathbf{M} , for any $m \times n$ matrix \mathbf{U} with orthonormal columns, a norm $\|\cdot\|_X$ is defined to be unitarily invariant if $\|\mathbf{U}\mathbf{M}\|_X = \|\mathbf{M}\|_X$. The Schatten-p norm is unitarily invariant for all $p \ge 1$.

There exists a closed-form expression for the low-rank approximation problem under Schatten-*p* norms:

Fact 4.9 (Schatten-*p* Low-Rank Approximation). Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and an integer $k \in \mathbb{N}$,

$$\mathbf{A}_k = \arg\min_{rank(\mathbf{X}) \le k} \|\mathbf{A} - \mathbf{X}\|_{\mathcal{S}_p},$$

where \mathbf{A}_k is the truncated SVD of \mathbf{A} .

Fact 4.10 (Araki–Lieb–Thirring Inequality [Ara90]). *Given PSD matrices* \mathbf{A} , $\mathbf{B} \in \mathbb{R}^{d \times d}$, *for any* $r \geq 1$, *the following inequality holds:*

$$\operatorname{Tr}\left(\left(\mathbf{B}\mathbf{A}\mathbf{B}\right)^{r}\right) \leq \operatorname{Tr}\left(\mathbf{B}^{r}\mathbf{A}^{r}\mathbf{B}^{r}\right).$$

Further, for 0 < r < 1, the reverse holds

$$\operatorname{Tr}\left(\left(\mathbf{B}\mathbf{A}\mathbf{B}\right)^{r}\right) \geq \operatorname{Tr}\left(\mathbf{B}^{r}\mathbf{A}^{r}\mathbf{B}^{r}\right).$$

Fact 4.11 (Mahler's Orthogonal Operator Inequality, Theorem 1.7 in [Mah90]). *Given* $p \ge 2$, and matrices **P** and **Q** such that the row (column) span of **P** is orthogonal to the row (column) span of **Q**, the following inequality holds:

$$\|\mathbf{P}\|_{\mathcal{S}_v}^p + \|\mathbf{Q}\|_{\mathcal{S}_v}^p \leq \|\mathbf{P} + \mathbf{Q}\|_{\mathcal{S}_v}^p$$
.

Fact 4.12 (Hölder's Inequality for Schatten-p Norms, Corollary 4.2.6 [Bha13]). *Given matrices* $\mathbf{A}, \mathbf{B}^{\mathsf{T}} \in \mathbb{R}^{n \times d}$ and $p \in [1, \infty)$, the following holds

$$\|\mathbf{A}\mathbf{B}\|_{\mathcal{S}_p} \leq \|\mathbf{A}\|_{\mathcal{S}_q} \cdot \|\mathbf{B}\|_{\mathcal{S}_r}$$
,

for any q, r such that $\frac{1}{p} = \frac{1}{q} + \frac{1}{r}$.

We also require *pinching inequalities* that were originally introduced to relate norms for partitioned operators over direct sums of Hilbert spaces. In our context, these inequalities simplify to norm inequalities for block matrices:

Fact 4.13 (Pinching Inequalities for Schatten-p Norms, [BKL02]). Let $\mathbf{M} \in \mathbb{R}^{td \times td}$ be the following block matrix

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_{(1,1)} & \mathbf{M}_{(1,2)} & \cdots & \mathbf{M}_{(1,t)} \\ \mathbf{M}_{(2,1)} & \mathbf{M}_{(2,2)} & \cdots & \mathbf{M}_{(1,t)} \\ \vdots & & \ddots & \vdots \\ \mathbf{M}_{(t,1)} & \mathbf{M}_{(t,2)} & \cdots & \mathbf{M}_{(t,t)} \end{bmatrix},$$

where for all $i, j \in [t]$, $\mathbf{M}_{(i,j)} \in \mathbb{R}^{d \times d}$. For all $p \ge 1$, the following inequality holds:

$$\left(\sum_{i\in[t]}\left\|\mathbf{M}_{(i,i)}\right\|_{\mathcal{S}_p}^p\right)^{1/p}\leq \|\mathbf{M}\|_{\mathcal{S}_p}.$$

We also require a norm compression inequality that is a special case of Conjecture 2.2 (and known to be true), when each block is aligned in the following sense:

Fact 4.14 (Aligned Norm Compression Inequality, Section 4.3 in [Aud08]). Let $\mathbf{M} = \begin{pmatrix} \mathbf{M}_1 & \mathbf{M}_2 \\ \mathbf{M}_3 & \mathbf{M}_4 \end{pmatrix}$ such that there exist scalars α_1 , α_2 , β_1 , β_2 such that $\mathbf{M}_1 = \alpha_1 \mathbf{X}$, $\mathbf{M}_2 = \alpha_2 \mathbf{X}$, $\mathbf{M}_3 = \beta_1 \mathbf{Y}$ and $\mathbf{M}_4 = \beta_2 \mathbf{Y}$. Then, for any $p \geq 2$,

$$\|\mathbf{M}\|_{\mathcal{S}_p} \leq \left\| \begin{pmatrix} \|\mathbf{M}_1\|_{\mathcal{S}_p} & \|\mathbf{M}_2\|_{\mathcal{S}_p} \\ \|\mathbf{M}_3\|_{\mathcal{S}_p} & \|\mathbf{M}_4\|_{\mathcal{S}_p} \end{pmatrix} \right\|_{\mathcal{S}_n}.$$

Random Matrix Theory. Next, we recall some basic facts for Wishart ensembles from random matrix theory (we refer the reader to [Tao12] for a comprehensive overview).

Definition 4.15 (Wishart Ensemble). An $n \times n$ matrix **W** is sampled from a Wishart Ensemble, Wishart(n), if **W** = $\mathbf{X}\mathbf{X}^{\mathsf{T}}$ such that for all $i, j \in [n] \mathbf{X}_{i,j} \sim \mathcal{N}\left(0, \frac{1}{n}\mathbf{I}\right)$.

Fact 4.16 (Norms of a Wishart Ensemble). Let $\mathbf{W} \sim \textit{Wishart}(n)$ such that $n = \Omega(1/\varepsilon^3)$. Then, with probability 99/100, $\|\mathbf{W}\|_{op} \leq 5$ and for any fixed constant p, $\|\mathbf{I} - \frac{1}{5}\mathbf{W}\|_{\mathcal{S}_v}^p = \Theta\left(\frac{1}{\varepsilon^{1/3}}\right)$.

5 Algorithms for Schatten-p LRA

In this section, we focus on obtaining algorithms for low-rank approximation in Schatten-p norm, simultaneously for all real, not necessarily constant, $p \in [1, O(\log(d)/\varepsilon)]$. For the special case of $p \in \{2, \infty\}$, Musco and Musco [MM15] showed an algorithm with matrix-vector query complexity $\tilde{O}(k/\varepsilon^{1/2})$, given below as Algorithm 5.6. We show that the number of matrix-vector products we require scales proportional to $\tilde{O}(kp^{1/6}/\varepsilon^{1/3})$ instead. Finally, recall when $p > \log(d)/\varepsilon$, it suffices to run Block Krylov for $p = \infty$, which requires $O(\log(d/\varepsilon)k/\sqrt{\varepsilon})$ matrix-vector products.

Theorem 5.1 (Optimal Schatten-p Low-Rank Approximation). Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, a target rank $k \in [d]$, an accuracy parameter $\epsilon \in (0,1)$ and any $p \in [1,O(\log(d)/\epsilon)]$, Algorithm 5.3 performs $O\left(\frac{kp^{1/6}\log(d/\epsilon)}{\epsilon^{1/3}} + \log(d/\epsilon)k\sqrt{p}\right)$ matrix-vector products and outputs a $d \times k$ matrix \mathbf{Z} with orthonormal columns such that with probability at least 9/10,

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p} \leq \left(1 + \epsilon \right) \min_{\mathbf{U}: \ \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{\mathcal{S}_p}.$$

Further, in the RAM model, the algorithm runs in time $O\left(\frac{\operatorname{nnz}(\mathbf{A})p^{1/6}k\log^2(d/\epsilon)}{\epsilon^{1/3}} + \frac{np^{(\omega-1)/6}k^{\omega-1}}{\epsilon^{(\omega-1)/3}}\right)$.

We first introduce the following lemmas from Musco and Musco [MM15] that provide convergence bounds for the performance of Block Krylov Iteration (Algorithm 5.6):

Lemma 5.2 (Gap Independent Block Krylov with Arbitrary Accuracy). Let **A** be an $n \times d$ matrix, k be the target rank and $\gamma > 0$ be an accuracy parameter. Then, initializing Algorithm 5.6 with block size k and running for $q = \Omega\left(\log(d/\gamma)/\sqrt{\gamma}\right)$ iterations outputs a $d \times k$ matrix **Z** such that with probability 99/100, for all $i \in [k]$,

$$\|\mathbf{A}\mathbf{Z}_{*,i}\|_{2}^{2} = \sigma_{i}^{2} \pm \gamma \sigma_{k+1}^{2}.$$

Further, the total number of matrix-vector products is O(kq) and the running time in the RAM model is $O(\operatorname{nnz}(\mathbf{A})kq + n(kq)^2 + (kq)^{\omega})$.

The aforementioned lemma follows directly from Theorem 1 in [MM15], using the per-vector error guarantee (3).

Algorithm 5.3 (Optimal Schatten-p Low-rank Approximation).

Input: An $n \times d$ matrix **A**, target rank $k \le d$, accuracy parameter $0 < \varepsilon < 1$, and $p \ge 1$.

- 1. Let $\gamma_1 = \varepsilon^{2/3}/p^{1/3}$. Run Block Krylov Iteration (Algorithm 5.6) on **A** with block size k, and number of iterations $q = O(\log(d/\gamma_1)/\sqrt{\gamma_1} + \log(d/\epsilon)\sqrt{p})$. Let $\mathbf{Z}_1 \in \mathbb{R}^{d \times k}$ be the corresponding output with orthonormal columns.
- 2. Run Block Krylov Iteration (Algorithm 5.6) on \mathbf{A}^{\top} with block size k, and number of iterations $q = O(\log(d/\gamma_1)/\sqrt{\gamma_1})$. Let $\mathbf{W}_1 \in \mathbb{R}^{n \times k}$ be the corresponding output with orthonormal columns.
- 3. Let $\gamma_2 = \varepsilon$ and let $s = O(p^{-1/3}k/\varepsilon^{1/3})$. Run Block Krylov Iteration (Algorithm 5.6) on \mathbf{A}^{\top} with block size s, and number of iterations $q = O(\log(d/\gamma_2)\sqrt{p})$. Let $\mathbf{W}_2 \in \mathbb{R}^{n \times k}$ be the corresponding output with orthonormal columns.
- 4. Run Block Krylov on **A** with target rank k+1 and number of iterations $q=O((\log(dp)+\log(d/\epsilon))\sqrt{p})$, and let $\hat{\mathbf{Z}}_1$ be the resulting $d\times(k+1)$ output matrix. Compute $\hat{\sigma}_1^2=\left\|\mathbf{A}(\hat{\mathbf{Z}}_1)_{*,1}\right\|_2^2$ and $\hat{\sigma}_{k+1}^2=\left\|\mathbf{A}(\hat{\mathbf{Z}}_1)_{*,k+1}\right\|_2^2$, rough estimates of the 1-st and (k+1)-st singular values of **A**. Run Block Krylov on **A** with target rank s, where $s=O(p^{-1/3}k/\epsilon^{1/3})$ and iterations $q=O(\log(d/\epsilon)\sqrt{p})$, and let $\hat{\mathbf{Z}}_2$ be the resulting $d\times s$ output matrix. Compute $\hat{\sigma}_s^2=\left\|\mathbf{A}(\hat{\mathbf{Z}}_2)_{*,s}\right\|_2^2$, an estimate to the s-th singular value of **A**.
- 5. If $\hat{\sigma}_1^2 \geq (1 + 0.5/p)\hat{\sigma}_{k+1}^2$, set $\mathbf{Z} = \mathbf{Z}_1$. Else, if $\hat{\sigma}_s^2 \leq \hat{\sigma}_{k+1}^2/(1 + 0.5/p)$, set \mathbf{Z} to be an orthonormal basis for $\mathbf{A}^{\top}\mathbf{W}_2\mathbf{W}_2^{\top}$ and otherwise set \mathbf{Z} to be an orthonormal basis for $\mathbf{A}^{\top}\mathbf{W}_1\mathbf{W}_1^{\top}$.

Output: A matrix $\mathbf{Z} \in \mathbb{R}^{d \times k}$ with orthonormal columns such that

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p}^p \leq (1 + \epsilon) \min_{\mathbf{U}: \ \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{\mathcal{S}_p}^p.$$

Lemma 5.4 (Gap Dependent Block Krylov, Theorem 13 [MM15]). Let **A** be an $n \times d$ matrix and $\gamma > 0$, be an accuracy parameter and $p, k \in N$ be such that $b \ge k$. Let $\sigma_1, \sigma_2 \dots \sigma_d$ be the singular values of **A**. Then, initializing Algorithm 5.6 with block size b and running for $q = \Omega\left(\log(n/\gamma)\sqrt{\sigma_k}/\sqrt{\sigma_k - \sigma_b}\right)$ iterations outputs a $d \times k$ matrix **Z** such that with probability 99/100, for all $i \in [k]$

$$\|\mathbf{A}\mathbf{Z}_{*,i}\|_{2}^{2} = \sigma_{i}^{2} \pm \gamma \sigma_{k+1}^{2}.$$

Further, the total number of matrix-vector products is O(pq) and the running time in the RAM model is $O(\operatorname{nnz}(\mathbf{A})bq + n(bq)^2 + (bq)^{\omega})$.

Next, we prove the following key lemma relating the Schatten-*p* norm of row and column projections applied to a matrix **A** to the Schatten-*p* norm of the matrix itself. We can interpret this lemma as an extension of the Pythagorean Theorem to Schatten-*p* spaces and believe this lemma is of independent interest. We note that we appeal to *pinching inequality* for partitioned operators to obtain this lemma.

Lemma 5.5 (Schatten-p Norms for Orthogonal Projections). Let **A** be an $n \times d$ matrix, let **P** be an $n \times n$ matrix, and let **Q** be a $d \times d$ matrix such that both **P** and **Q** are orthogonal projection matrices of rank k (see Definition 4.1). Then, the following inequality holds for all $p \ge 1$:

$$\|\mathbf{A}\|_{\mathcal{S}_{v}}^{p} \geq \|\mathbf{P}\mathbf{A}\mathbf{Q}\|_{\mathcal{S}_{v}}^{p} + \|(\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_{v}}^{p}.$$

Proof. Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}$ be the SVD of \mathbf{A} , where $\mathbf{U} \in \mathbb{R}^{n \times d}$ and $\mathbf{V}^{\top} \in \mathbb{R}^{d \times d}$ have orthonormal columns and rows respectively. We construct unitary matrices \mathbf{R} and \mathbf{S} , such that $\mathbf{R} \in \mathbb{R}^{n \times n}$ and $\mathbf{S} \in \mathbb{R}^{d \times d}$ that satisfy the following constraints:

1. $\mathbf{R}^{\mathsf{T}}\mathbf{I}_{k}\mathbf{R}\mathbf{A}\mathbf{S}^{\mathsf{T}}\mathbf{I}_{k}\mathbf{S} = \mathbf{P}\mathbf{A}\mathbf{Q}$, and

2.
$$\mathbf{R}^{\top} (\mathbf{I} - \mathbf{I}_k) \mathbf{R} \mathbf{A} \mathbf{S}^{\top} (\mathbf{I} - \mathbf{I}_k) \mathbf{S} = (\mathbf{I} - \mathbf{P}) \mathbf{A} (\mathbf{I} - \mathbf{Q}),$$

where the trunctated Identity matrix, I_k , left multiplying A is $n \times n$ and right multiplying A is $d \times d$. Recall, since P is a rank-k projection matrix, it admits a decomposition $P = XX^T$ such that X has k orthonormal columns and similarly $I - P = YY^T$, where Y has n - k orthonormal columns. Further, since X and Y span disjoint subspaces, and the union of their span is \mathbb{R}^n , the matrix $(X \mid Y)$, obtained by concatenating their columns, is unitary. Then, it suffices to set $R = (X \mid Y)^T$. To see this, observe,

$$\mathbf{R}^{\mathsf{T}}\mathbf{I}_{k}\mathbf{R} = (\mathbf{X} \mid 0) \cdot \begin{pmatrix} \mathbf{X}^{\mathsf{T}} \\ 0 \end{pmatrix} = \mathbf{X}\mathbf{X}^{\mathsf{T}} = \mathbf{P},$$

and similarly,

$$\mathbf{R}^{\top} \left(\mathbf{I} - \mathbf{I}_{k} \right) \mathbf{R} = \mathbf{Y} \mathbf{Y}^{\top} = \mathbf{I} - \mathbf{P}.$$

We repeat the above argument for the projection matrix \mathbf{Q} . Let $\mathbf{Q} = \mathbf{W}\mathbf{W}^{\mathsf{T}}$, where \mathbf{W} is $d \times k$ and has orthonormal columns, and $\mathbf{I} - \mathbf{Q} = \mathbf{Z}\mathbf{Z}^{\mathsf{T}}$, where \mathbf{Z} is $d \times (d - k)$ and has orthonormal columns. Observe, it suffices to set $\mathbf{S} = (\mathbf{W} \mid \mathbf{Z})^{\mathsf{T}}$, since \mathbf{S} is unitary and $\mathbf{S}^{\mathsf{T}}\mathbf{I}_k\mathbf{S} = \mathbf{Q}$ and $\mathbf{S}^{\mathsf{T}}(\mathbf{I} - \mathbf{I}_k)\mathbf{S} = \mathbf{I} - \mathbf{Q}$. Note, by construction, we satisfy the two aforementioned constraints.

Let $\hat{\mathbf{A}} = \mathbf{R}\mathbf{A}\mathbf{S}^{\mathsf{T}}$. Since \mathbf{R} and \mathbf{S} are unitary, it follows from unitary invariance of the Schatten-p norm that

$$\left\| \hat{\mathbf{A}} \right\|_{\mathcal{S}_{v}} = \left\| \mathbf{R} \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\mathsf{T}} \mathbf{S}^{\mathsf{T}} \right\|_{\mathcal{S}_{p}} = \left\| \mathbf{A} \right\|_{\mathcal{S}_{p}}$$
 (5.1)

Further, observe for any $n \times d$ matrix **M**, we have the following block decomposition

$$\mathbf{M} = \mathbf{I}_{k} \mathbf{M} \mathbf{I}_{k} + \mathbf{I}_{k} \mathbf{M} \left(\mathbf{I} - \mathbf{I}_{k} \right) + \left(\mathbf{I} - \mathbf{I}_{k} \right) \mathbf{M} \mathbf{I}_{k} + \left(\mathbf{I} - \mathbf{I}_{k} \right) \mathbf{M} \left(\mathbf{I} - \mathbf{I}_{k} \right)$$

$$= \begin{pmatrix} \mathbf{M}_{1:k,1:k} & \mathbf{M}_{1:k,k+1:d} \\ \mathbf{M}_{k+1:n,1:k} & \mathbf{M}_{k+1:n,k+1:d} \end{pmatrix},$$

where the notation $\mathbf{M}_{i:i',j:j'}$ picks the $(i'-i+1)\times(j'-j+1)$ sized sub-matrix corresponding to the rows indices [i, i'] and column indices [j, j']. Since appending rows and columns of 0's does not change the singular values, we have $\|\mathbf{I}_{k}\mathbf{M}\mathbf{I}_{k}\|_{\mathcal{S}_{p}} = \|\mathbf{M}_{1:k,1:k}\|_{\mathcal{S}_{p}}$ and $\|(\mathbf{I} - \mathbf{I}_{k})\mathbf{M}(\mathbf{I} - \mathbf{I}_{k})\|_{\mathcal{S}_{p}} = \|\mathbf{M}_{k+1:n,k+1:d}\|_{\mathcal{S}_{p}}$. Setting $\mathbf{M} = \hat{\mathbf{A}}$, we have

$$\begin{aligned} \left\| \hat{\mathbf{A}} \right\|_{\mathcal{S}_{p}}^{p} &= \left\| \begin{pmatrix} \hat{\mathbf{A}}_{1:k,1:k} & \hat{\mathbf{A}}_{1:k,k+1:d} \\ \hat{\mathbf{A}}_{k+1:n,1:k} & \hat{\mathbf{A}}_{k+1:n,k+1:d} \end{pmatrix} \right\|_{\mathcal{S}_{p}}^{p} \\ &\geq \left\| \hat{\mathbf{A}}_{1:k,1:k} \right\|_{\mathcal{S}_{p}}^{p} + \left\| \hat{\mathbf{A}}_{k+1:n,k+1:d} \right\|_{\mathcal{S}_{p}}^{p} \\ &= \left\| \mathbf{I}_{k} \hat{\mathbf{A}} \mathbf{I}_{k} \right\|_{\mathcal{S}_{p}}^{p} + \left\| (\mathbf{I} - \mathbf{I}_{k}) \hat{\mathbf{A}} (\mathbf{I} - \mathbf{I}_{k}) \right\|_{\mathcal{S}_{p}}^{p} , \end{aligned}$$

$$(5.2)$$

where the inequality follows from using the *pinching inequality* on the block matrix (see Fact 4.13). By the unitary invariance of the Schatten-*p* norm, we have

$$\left\|\mathbf{I}_{k}\mathbf{\hat{A}}\mathbf{I}_{k}\right\|_{\mathcal{S}_{p}}^{p} = \left\|\mathbf{R}^{\top}\mathbf{I}_{k}\mathbf{\hat{A}}\mathbf{I}_{k}\mathbf{S}\right\|_{\mathcal{S}_{p}}^{p} = \left\|\mathbf{P}\mathbf{A}\mathbf{Q}\right\|_{\mathcal{S}_{p}}^{p},$$

and similarly,

$$\left\| \left(\mathbf{I} - \mathbf{I}_k \right) \hat{\mathbf{A}} \left(\mathbf{I} - \mathbf{I}_k \right) \right\|_{\mathcal{S}_v}^p = \left\| \mathbf{R}^\top \left(\mathbf{I} - \mathbf{I}_k \right) \hat{\mathbf{A}} \left(\mathbf{I} - \mathbf{I}_k \right) \mathbf{S} \right\|_{\mathcal{S}_v}^p = \left\| \left(\mathbf{I} - \mathbf{P} \right) \mathbf{A} \left(\mathbf{I} - \mathbf{Q} \right) \right\|_{\mathcal{S}_p}^p.$$

Plugging these two bounds back into Equation (5.2), along with Equation (5.1), we can conclude,

$$\left\|\mathbf{A}\right\|_{\mathcal{S}_{p}}^{p} \geq \left\|\mathbf{P}\mathbf{A}\mathbf{Q}\right\|_{\mathcal{S}_{p}}^{p} + \left\|\left(\mathbf{I} - \mathbf{P}\right)\mathbf{A}\left(\mathbf{I} - \mathbf{Q}\right)\right\|_{\mathcal{S}_{p}}^{p}.$$

Algorithm 5.6 (Block Krylov Iteration, [MM15]).

Input: An $n \times d$ matrix **A**, target rank k, iteration count q and a block size parameter s such that $k \le s \le d$.

- 1. Let **U** be a $n \times s$ matrix such that each entry is drawn i.i.d. from $\mathcal{N}(0,1)$. Let $\mathbf{K} = \left[\mathbf{A}^{\mathsf{T}} \mathbf{U}; (\mathbf{A}^{\mathsf{T}} \mathbf{A}) \mathbf{A}^{\mathsf{T}} \mathbf{U}; (\mathbf{A}^{\mathsf{T}} \mathbf{A})^2 \mathbf{A}^{\mathsf{T}} \mathbf{U}; \dots; (\mathbf{A}^{\mathsf{T}} \mathbf{A})^q \mathbf{A}^{\mathsf{T}} \mathbf{U} \right]$ be the $d \times s(q+1)$ Krylov matrix obtained by concatenating the matrices $\mathbf{A}^{\mathsf{T}}\mathbf{U}, \dots, (\mathbf{A}^{\mathsf{T}}\mathbf{A})^q \mathbf{A}^{\mathsf{T}}\mathbf{U}$.
- 2. Compute an orthonomal basis **Q** for the column span of **K**. Let $\mathbf{M} = \mathbf{Q}^{\mathsf{T}} \mathbf{A}^{\mathsf{T}} \mathbf{A} \mathbf{Q}$.
- 3. Compute the top k left singular vectors of \mathbf{M} , and denote them by \mathbf{Y}_k .

Output: $Z = QY_k$

Note, despite establishing Lemma 5.5, it is not immediately apparent how to lower bound $\|\mathbf{AZZ}^{\mathsf{T}}\|_{\mathbf{S}_n}^p$, where **Z** is a candidate solution. Next, we show how to translate a guarantee on the Euclidean norm of **A** times a column of **Z** to a lower bound on $\|\mathbf{A}\mathbf{Z}\mathbf{Z}^{\mathsf{T}}\|_{\mathcal{S}_{\bullet}}^{p}$.

Lemma 5.7 (Per-Vector Guarantees to Schatten Norms). Let **A** be an $n \times d$ matrix with singular values denoted by $\{\sigma_i(\mathbf{A})\}_{i \in [d]}$. Let **Z** be a $d \times k$ matrix with orthonormal columns that is output by Algorithm 5.6, such that for all $i \in [k]$, with probability at least 99/100, $\|\mathbf{AZ}_{*,i}\|_2^2 \ge \sigma_i^2(\mathbf{A}) - \gamma \sigma_{k+1}^2(\mathbf{A})$, for some $\gamma \in (0,1)$. Then, for any $p \ge 1$, we have

$$\|\mathbf{A}\mathbf{Z}\mathbf{Z}^{\top}\|_{\mathcal{S}_{p}}^{p} \geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - O(\gamma p) \sum_{i \in [k]} \sigma_{k+1}^{2}(\mathbf{A}) \sigma_{i}^{p-2}(\mathbf{A}).$$

Proof. First, we observe that it suffices to show that $\sigma_i(\mathbf{AZ})^2 \ge \|\mathbf{A}z_i\|_2^2$, where z_i is shorthand for $\mathbf{Z}_{*,i}$, the *i*-th column of \mathbf{Z} . Assuming this inequality holds, we can complete the proof as follows: we know that for all $i \in [k]$,

$$\sigma_i^2(\mathbf{AZ}) \ge \|\mathbf{A}z_i\|_2^2 \ge \sigma_i^2(\mathbf{A}) - \gamma \sigma_{k+1}^2(\mathbf{A})$$

$$= \sigma_i^2(\mathbf{A}) \left(1 - \gamma \frac{\sigma_{k+1}^2(\mathbf{A})}{\sigma_i^2(\mathbf{A})} \right)$$
(5.3)

Then, taking p/2-th powers in (5.3),

$$\sigma_{i}^{p}(\mathbf{AZ}) \geq \sigma_{i}^{p}(\mathbf{A}) \left(1 - \gamma \frac{\sigma_{k+1}^{2}(\mathbf{A})}{\sigma_{i}^{2}(\mathbf{A})}\right)^{p/2}$$

$$\geq \sigma_{i}^{p}(\mathbf{A}) \left(1 - O\left(\frac{\gamma p \sigma_{k+1}^{2}(\mathbf{A})}{\sigma_{i}^{2}(\mathbf{A})}\right)\right)$$

$$= \sigma_{i}^{p}(\mathbf{A}) - O(\gamma p) \sigma_{k+1}^{2}(\mathbf{A}) \sigma_{i}^{p-2}(\mathbf{A})$$
(5.4)

where the second inequality follows from the generalized Bernoulli inequality (see Fact 4.6). Summing over all $i \in [k]$, we can conclude

$$\|\mathbf{A}\mathbf{Z}\|_{\mathcal{S}_{p}}^{p} \geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - \sum_{i \in [k]} O(\gamma p) \sigma_{k+1}^{2}(\mathbf{A}) \sigma_{i}^{p-2}(\mathbf{A}).$$

Therefore, it remains to show that $\sigma_i(\mathbf{AZ})^2 \ge \|\mathbf{A}z_i\|_2^2$. First, we recall that Algorithm 5.6 outputs $\{z_i\}_{i\in[k]}$ such that $z_i = \mathbf{Q}\tilde{z}_i$, where \mathbf{Q} is an orthonormal basis for the Krylov space \mathbf{K} (an $d\times s(q+1)$ matrix) and \tilde{z}_i is the i-th singular vector of $\mathbf{Q}^{\top}\mathbf{A}^{\top}\mathbf{AQ}$. Note that the \tilde{z}_i 's are s(q+1)-dimensional vectors. Let $\mathbf{W}\Omega\mathbf{W}^{\top}$ be the SVD of $\mathbf{Q}^{\top}\mathbf{A}^{\top}\mathbf{AQ}$. Then, $\mathbf{Q}\mathbf{W}\Omega\mathbf{W}^{\top}\mathbf{Q}^{\top}$ is the SVD of $\mathbf{Q}\mathbf{Q}^{\top}\mathbf{A}^{\top}\mathbf{AQ}\mathbf{Q}^{\top}$. To see this, let the i-th column of $\mathbf{Q}\mathbf{W}$ be denoted by $\mathbf{Q}\mathbf{W}_{*,i}$. Then,

$$\langle \mathbf{Q}\mathbf{W}_{*,i}, \mathbf{Q}\mathbf{W}_{*,i} \rangle = \mathbf{W}_{*,i}^{\mathsf{T}} \mathbf{Q}^{\mathsf{T}} \mathbf{Q} \mathbf{W}_{*,i} = 1$$

and similarly for any $j \neq i$,

$$\left\langle \mathbf{Q}\mathbf{W}_{*,i},\mathbf{Q}\mathbf{W}_{*,j}\right\rangle = \mathbf{W}_{*,i}^{\top}\mathbf{Q}^{\top}\mathbf{Q}\mathbf{W}_{*,j} = 0$$

where we use that $\mathbf{Q}^{\mathsf{T}}\mathbf{Q} = \mathbf{I}$ and the columns of \mathbf{W} are orthonormal, which holds by definition. Therefore, $z_i = \mathbf{Q}\tilde{z}_i$ is the *i*-th singular vector of $\mathbf{Q}\mathbf{Q}^{\mathsf{T}}\mathbf{A}^{\mathsf{T}}\mathbf{A}\mathbf{Q}\mathbf{Q}^{\mathsf{T}}$. Let $\tilde{\mathbf{Z}}$ be the matrix obtained by stacking the vectors \tilde{z}_i together. Then, we have

$$\sigma_{i}(\mathbf{A}\mathbf{Z})^{2} = \sigma_{i}^{2}(\mathbf{A}\mathbf{Q}\tilde{\mathbf{Z}}) = \sigma_{i}^{2}(\mathbf{A}\mathbf{Q})$$

$$= \sigma_{i}^{2}(\mathbf{A}\mathbf{Q}\mathbf{Q}^{\mathsf{T}})$$

$$= z_{i}^{\mathsf{T}}\mathbf{Q}\mathbf{Q}^{\mathsf{T}}\mathbf{A}^{\mathsf{T}}\mathbf{A}\mathbf{Q}\mathbf{Q}^{\mathsf{T}}z_{i}$$

$$= z_{i}^{\mathsf{T}}\mathbf{A}^{\mathsf{T}}\mathbf{A}z_{i}$$
(5.5)

where the first equality follows from the definition of $\tilde{\mathbf{Z}}$, the second follows from observing that $\tilde{\mathbf{Z}}$ are the singular vectors of $\mathbf{A}\mathbf{Q}$ as shown above, the third follows from \mathbf{Q}^{T} having orthonormal rows, the fourth from z_i being the i-th singular vector of $\mathbf{A}\mathbf{Q}\mathbf{Q}^{\mathsf{T}}$ and the last from observing that z_i is in the column span of \mathbf{Q} and thus $\mathbf{Q}\mathbf{Q}^{\mathsf{T}}z_i=z_i$. This concludes the proof.

Finally, we also need the following lemma:

Lemma 5.8 (Singular Values to Alignment of Singular Vectors). Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}$ be the SVD and let \mathbf{Z} be a $d \times k$ orthonormal matrix such that for all $i \in [k]$, $\|\mathbf{A} \mathbf{Z}_{*,i}\|_2^2 \ge \sigma_i^2(\mathbf{A}) - (\epsilon/d)^c \sigma_{k+1}^2$, for some fixed constant $c \ge 10$. Further, assume there exists a $j^* \in [k]$ such that for all $j \in [j^*]$, $\sigma_j^2(\mathbf{A}) \ge (1 + \epsilon/d) \sigma_{k+1}^2(\mathbf{A})$ and $\sigma_{j^*+1}^2(\mathbf{A}) \le (1 - \epsilon/d) \sigma_{j^*}^2(\mathbf{A})$. Then,

$$\left\|\mathbf{V}_{j^*}^{\top}\mathbf{Z}\right\|_F^2 \ge j^* - (\epsilon/d)^{c-4},$$

where $\mathbf{V}_{j^*}^{\mathsf{T}}$ is the top- j^* rows of \mathbf{V}^{T} .

Proof. First, using our hypothesis and summing over all $\ell \in [j^*]$, we have

$$\sum_{\ell \in [j^*]} \|\mathbf{A} \mathbf{Z}_{*,\ell}\|_2^2 \ge \|\mathbf{A}_{j^*}\|_F^2 - (\epsilon/d)^{c-1} \sigma_{k+1}^2.$$
 (5.6)

Let **B** be a $d \times k$ matrix with entries $b_{j\ell} = (\mathbf{V}_{j,*}^{\top} \mathbf{Z}_{*,\ell})^2$, and let $v_j = \sum_{\ell \in [j^*]} b_{j\ell}$. Using this notation, and since **V** and **Z** are orthonormal we have

$$v_j \le \sum_{\ell \in [k]} b_{j\ell} \le 1 \text{ for } j \in [d]$$

$$(5.7)$$

$$\sum_{j \in [d]} b_{j\ell} \le 1 \text{ for } \ell \in [k]$$
(5.8)

$$\sum_{\ell \in [j^*]} \|\mathbf{A} \mathbf{Z}_{*,\ell}\|_2^2 = \sum_{\substack{j \in [d] \\ \ell \in [j^*]}} \sigma_j^2 b_{j\ell} = \sum_{j \in [d]} \sigma_j^2 v_j, \tag{5.9}$$

where we abbreviate $\sigma_i(\mathbf{A})$ by σ_i in this proof. Define α by

$$\alpha \equiv j^* - \left\| \mathbf{V}_{j^*}^\top \mathbf{Z} \right\|_F^2 = j^* - \sum_{\substack{j \in [j^*] \\ \ell \in [k]}} b_{j\ell} \le j^* - \sum_{j \in [j^*]} v_j, \tag{5.10}$$

so the claim of the lemma is an upper bound on α . If $\alpha \le 0$ the lemma follows, so assume $\alpha > 0$. We will show the inequality

$$\sum_{\ell \in [j^*]} \|\mathbf{A} \mathbf{Z}_{*,\ell}\|_2^2 = \sum_{j \in [d]} \sigma_j^2 v_j \le \|\mathbf{A}_{j^*}\|_F^2 - \left(\frac{\epsilon \alpha}{d}\right) \sigma_{j^*}^2.$$
 (5.11)

Assuming this, and comparing it with the lower bound on the LHS in (5.6), we can then conclude that

$$\left(\frac{\epsilon}{d}\right)^{c-1}\sigma_{k+1}^2 \ge \frac{\epsilon\alpha}{d}\sigma_{j^*}^2 \tag{5.12}$$

which in turn bounds $\alpha \leq \left(\frac{\epsilon}{d}\right)^{c-3}$, which is the claim of the lemma.

It remains to show (5.11). We have constraints (5.7) and (5.10) on v_i , and also

$$\sum_{j \in [d]} v_j = \sum_{\ell \in [j^*]} \sum_{j \in [d]} b_{j\ell} = \sum_{\ell \in [j^*]} \sum_{j \in [d]} b_{j\ell} \le j^* \text{ by (5.8)}$$
(5.13)

The maximum value of $\sum_{j \in [d]} \sigma_j^2 v_j$ under constraints (5.7), (5.10), and (5.13) results by pushing all the "weight" of $\sum_j v_j$ to the larger σ_j^2 to the maximum extent possible, that is, for $\hat{j} \equiv \lfloor j^* - \alpha \rfloor$, setting

$$\begin{aligned} v_j &\leftarrow 1 \text{ for } j \in [\hat{j}] \\ v_{\hat{j}+1} &\leftarrow j^* - \alpha - \hat{j} \\ v_j &\leftarrow 1 \text{ for } j = j^* + 1, \dots, j^* + \lfloor \alpha \rfloor \\ v_{j^* + \lfloor \alpha \rfloor + 1} &\leftarrow \alpha - \lfloor \alpha \rfloor. \end{aligned}$$

This is under the assumption that $\sum_{j \in [d]} v_j$ is equal to its upper bound; it might be smaller, but if so, $\sum_{j \in [d]} \sigma_j^2 v_j$ can only be smaller. However, if $\alpha \ge 1$, then $v_{j*+1} = 1$ in the above, and by hypothesis $\sigma_{j^*+1}^2 \le (1 - \epsilon/d) \sigma_{j^*}^2$, and so $\sum_j \sigma_j^2 v_j \le \|A_{j^*}\|_F^2 - (\epsilon/d) \sigma_{j^*}^2$, contradicting (5.6). So $\alpha < 1$, and the above simplifies to

$$v_j \leftarrow 1 \text{ for } j \in [j^* - 1]$$

 $v_{j^*} \leftarrow 1 - \alpha$
 $v_{j^*+1} \leftarrow \alpha$.

With this maximizing assignment, we have:

$$\left\| \mathbf{V}_{j^*}^{\top} \mathbf{Z} \right\|_F^2 = \sum_{j \in [d]} \sigma_j^2 v_j \le \sum_{j \in [j^*-1]} \sigma_j^2 + (1-\alpha) \sigma_{j^*}^2 + \alpha \sigma_{j^*+1}^2 \le \left\| \mathbf{A}_{j^*} \right\|_F^2 - \left(\frac{\epsilon \alpha}{d} \right) \sigma_{j^*}^2,$$

proving (5.11), which then implies the lemma as discussed.

Finally, we need a lemma relating the Schatten-p norm of **AZ** to that of $\mathbf{W}^{\mathsf{T}}\mathbf{A}$, where **Z** is an arbitrary orthonormal basis and **W** is an orthonormal basis for **AZ**.

Lemma 5.9. Given a full-rank $n \times d$ matrix \mathbf{A} , let \mathbf{Z} be a $d \times k$ matrix with orthonormal columns. Further, let \mathbf{W} be an $n \times k$ matrix with orthonormal columns such that \mathbf{W} is a basis for \mathbf{AZ} . Then, for all $i \in [k]$,

$$\sigma_i (\mathbf{W}^{\mathsf{T}} \mathbf{A})^p \geq \sigma_i (\mathbf{A} \mathbf{Z})^p$$

Proof. We use the following fact that for two matrices **A** and **B**, we have that for all i, $\sigma_i(\mathbf{A} \cdot \mathbf{B}) \leq \sigma_i(\mathbf{A}) \cdot \sigma_1(\mathbf{B})$; see, e.g., (2) in [LC15] and references [33-36] therein.

Using this fact, we have

$$\sigma_i(\mathbf{AZ}) = \sigma_i(\mathbf{AZZ}^T) = \sigma_i(\mathbf{WW}^T\mathbf{AZZ}^T) \le \sigma_i(\mathbf{WW}^T\mathbf{A}) \cdot \sigma_1(\mathbf{ZZ}^T) = \sigma_i(\mathbf{WW}^T\mathbf{A}) = \sigma_i(\mathbf{W}^T\mathbf{A}),$$

where we have used that $\sigma_1(\mathbf{Z}\mathbf{Z}^T) = 1$ since $\mathbf{Z}\mathbf{Z}^T$ is a projection matrix, and the fact that $\mathbf{W}\mathbf{W}^T$ is a basis for the column span of $\mathbf{A}\mathbf{Z}$. Raising both sides to the p-th power establishes the lemma.

We now have all the ingredients we need to complete the proof of Theorem 5.1.

Proof of Theorem 5.1. Observe, using Lemma 5.2 with probability at least 97/100, Step 3 of Algorithm 5.3 outputs $\hat{\sigma}_1^2 = (1 \pm 0.1/p) \, \sigma_1^2$, $\hat{\sigma}_{k+1}^2 = (1 \pm 0.1/p) \, \sigma_{k+1}^2$ and $\hat{\sigma}_s^2 = (1 \pm 0.1/p) \, \sigma_s^2$, for $s = O(kp^{-1/3}/\epsilon^{1/3})$. Condition on this event. Our proof proceeds via case analysis. The case where there is at least a constant gap between the first and (k+1)-st singular value is easy to handle since we can use gap-dependent guarantees to obtain highly accurate estimates of the top-k singular values. When there is no gap, either the Schatten-k norm of the tail is large compared to the k 1)-st singular value, and we don't require a highly accurate solution, or the Schatten-k norm of the tail is small, and increasing the block size induces a gap. We formalize this intuition into a proof.

Let us first consider the case where there is a constant gap between the top and the (k + 1)-st singular values, i.e., $\sigma_1 > (1 + 1/p)\sigma_{k+1}$. Observe, since we have (1 + 0.1/p)-approximate estimates to σ_1 and σ_{k+1} , we can detect that we are in this case and Algorithm 5.3 outputs $\mathbf{Z} = \mathbf{Z}_1$. We further observe that the Algorithm 5.3 runs at least $\Omega(\log(d/\epsilon)\sqrt{p})$ iterations (since $p \leq \log(d)/\epsilon$) since $\mathbf{Z} = \mathbf{Z}_1$. We observe that in this case, there exists a gap of size p between σ_1 and σ_{k+1} , since $1 - \sigma_{k+1}/\sigma_1 \leq 1/p$. It follows from Lemma 5.4 that running Block Krylov Iteration for $O(\log(d/\epsilon)\sqrt{p})$ iterations with block size $\geq k$ suffices to output a matrix \mathbf{Z} such that with probability at least 99/100, for all $i \in [k]$,

$$\|\mathbf{A}\mathbf{Z}_{*,i}\|_{2}^{2} \ge \sigma_{i}^{2}(\mathbf{A}) - \operatorname{poly}\left(\frac{\epsilon}{d}\right)\sigma_{k+1}^{2}(\mathbf{A}).$$
 (5.14)

We note that we cannot simply take p/2-th powers here (for large p) as this would introduce cross terms that scale proportional to $\sigma_i(\mathbf{A})$, which can be significantly larger than $\sigma_{k+1}(\mathbf{A})$. Instead, we require a finer analysis by splitting \mathbf{A} into a head and tail term.

Let $\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\mathsf{T}}$ be the SVD of \mathbf{A} and for all $j \in [d]$, let $v_j = \mathbf{V}_{j,*}^{\mathsf{T}}$ denote the j-th row of \mathbf{V}^{T} . By

the Pythagorean Theorem, we have

$$\|\mathbf{A}\mathbf{Z}\|_{F}^{2} = \|\mathbf{A}_{k}\mathbf{Z}\|_{F}^{2} + \|(\mathbf{A} - \mathbf{A}_{k})\mathbf{Z}\|_{F}^{2}$$

$$\leq \|\mathbf{\Sigma}_{k}\mathbf{V}_{k}^{\top}\mathbf{Z}\|_{F}^{2} + \sigma_{k+1}^{2} \|(\mathbf{I} - \mathbf{V}_{k}\mathbf{V}_{k}^{\top})\mathbf{Z}\|_{F}^{2}$$

$$= \sum_{j \in [k]} \sigma_{j}^{2} \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2} + \sigma_{k+1}^{2} (\|\mathbf{Z}\|_{F}^{2} - \sum_{j \in [k]} \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2})$$

$$= \sum_{j \in [k]} (\sigma_{j}^{2} - \sigma_{k+1}^{2}) \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2} + \sigma_{k+1}^{2}k.$$
(5.15)

Summing over $j \in [k]$ for the guarantee obtained in Equation 5.14, we have

$$\|\mathbf{A}\mathbf{Z}\|_{F}^{2} = \sum_{j \in [k]} \|\mathbf{A}\mathbf{Z}_{*,j}\|_{F}^{2} \ge \sum_{j \in [k]} \sigma_{j}^{2} - O(\gamma k) \sigma_{k+1}^{2}.$$
(5.16)

where $\gamma = \text{poly}(\epsilon/d)$. Combining Equations (5.15) and (5.16), we can conclude

$$\sum_{j \in [k]} \left(\sigma_j^2 - \sigma_{k+1}^2 \right) - O(\gamma k) \, \sigma_{k+1}^2 \le \sum_{j \in [k]} \left(\sigma_j^2 - \sigma_{k+1}^2 \right) \left\| v_j^\top \mathbf{Z} \right\|_2^2. \tag{5.17}$$

Let $j' \in [k]$ be the largest integer such that for all $j \le j'$, $\sigma_j^2 \ge (1 + \epsilon/d) \sigma_{k+1}^2$. Next, let $j^* \in [j', k]$ be such that $\sigma_{j^*+1} \le (1 - \epsilon/d)\sigma_{j^*}$. Observe, such a j^* is guaranteed to exist since there is a gap between σ_1 and σ_{k+1} . Since $\left\|v_j^\top \mathbf{Z}\right\|_2^2 \le 1$, we can restate Equation (5.17), as follows:

$$\begin{split} \sum_{j \in [k]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right) - O(\gamma k) \, \sigma_{k+1}^{2} &\leq \sum_{j \in [j^{*}]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right) \left\| v_{j}^{\mathsf{T}} \mathbf{Z} \right\|_{2}^{2} + \sum_{j \in [j^{*}+1,k]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right) \left\| v_{j}^{\mathsf{T}} \mathbf{Z} \right\|_{2}^{2} \\ &\leq \sum_{j \in [j^{*}]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right) \left\| v_{j}^{\mathsf{T}} \mathbf{Z} \right\|_{2}^{2} + \sum_{j \in [j^{*}+1,k]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right). \end{split}$$

Subtracting $\sum_{j \in [j^*+1,k]} \left(\sigma_j^2 - \sigma_{k+1}^2\right)$ from both sides, and rearranging, we have

$$\sum_{j \in [j^*]} \left(\sigma_j^2 - \sigma_{k+1}^2 \right) - O(\gamma k) \, \sigma_{k+1}^2 + \sigma_{k+1}^2 \sum_{j \in [j^*]} \left\| v_j^\top \mathbf{Z} \right\|_2^2 \le \sum_{j \in [j^*]} \sigma_j^2 \left\| v_j^\top \mathbf{Z} \right\|_2^2. \tag{5.18}$$

We are now ready to bound $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_{v}}$. By the triangle inequality,

$$\|\mathbf{A}\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\|_{\mathcal{S}_{p}} \leq \|\mathbf{A}_{j^{*}}\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\|_{\mathcal{S}_{p}} + \|\left(\mathbf{A} - \mathbf{A}_{j^{*}}\right)\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\|_{\mathcal{S}_{p}}$$
(5.19)

Observe, for any $p \ge 1$, $\|\mathbf{A}_{j^*}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\mathsf{T}})\|_{\mathcal{S}_p} \le \sqrt{k} \|\mathbf{A}_{j^*}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\mathsf{T}})\|_F$, since \mathbf{A}_{j^*} has rank at most k, with p = 1 achieving the worst inequality. Therefore, using the Pythagorean theorem again, and

plugging in the lower bound from Equation (5.18)

$$\|\mathbf{A}_{j^{*}} (\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_{p}} \leq \sqrt{k} \cdot \left(\sum_{j \in [j^{*}]} \sigma_{j}^{2} - \sigma_{j}^{2} \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2} \right)^{1/2}$$

$$\leq \sqrt{k} \cdot \left(\sum_{j \in [j^{*}]} \sigma_{j}^{2} - \left(\sum_{j \in [j^{*}]} \left(\sigma_{j}^{2} - \sigma_{k+1}^{2} \right) - O(\gamma k) \sigma_{k+1}^{2} + \sigma_{k+1}^{2} \sum_{j \in [j^{*}]} \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2} \right) \right)^{1/2}$$

$$\leq \sqrt{k} \sigma_{k+1} \cdot \left(j^{*} - \sum_{j \in [j^{*}]} \|v_{j}^{\top}\mathbf{Z}\|_{2}^{2} + O(\gamma k) \right)^{1/2}$$
(5.20)

It therefore remains to lower bound $\sum_{j \in [j^*]} \left\| v_j^{\mathsf{T}} \mathbf{Z} \right\|_2^2$. Applying Lemma 5.8, we have,

$$\sum_{j \in [j^*]} \left\| v_j^{\top} \mathbf{Z} \right\|_2^2 = \left\| \mathbf{V}_{j^*}^{\top} \mathbf{Z} \right\|_F^2 \ge j^* - O((\epsilon/d)^4)$$
 (5.21)

Plugging back into Equation (5.20), $\|\mathbf{A}_{j^*}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p} \leq O(\frac{\epsilon}{d}\sigma_{k+1})$ and thus substituting into Equation (5.19),

$$\left\|\mathbf{A}\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\right\|_{\mathcal{S}_{p}} \leq O\left(\frac{\epsilon}{d}\right) \left\|\mathbf{A} - \mathbf{A}_{k}\right\|_{\mathcal{S}_{p}} + \underbrace{\left\|\left(\mathbf{A} - \mathbf{A}_{j^{*}}\right)\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\right\|_{\mathcal{S}_{p}}}_{5.22.1}.$$
(5.22)

It remains to bound term 5.22.1 above.

Applying Lemma 5.5 with $\mathbf{Q} = \mathbf{Z}\mathbf{Z}^{\top}$ and $\mathbf{P} = \mathbf{W}\mathbf{W}^{\top}$ being the projection on the column span of $\mathbf{A}\mathbf{Z}\mathbf{Z}^{\top}$, we have

$$\begin{aligned} \left\| \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p}^p &\leq \left\| \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \right\|_{\mathcal{S}_p}^p - \left\| \mathbf{W} \mathbf{W}^{\top} \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \right\|_{\mathcal{S}_p}^p \\ &= \sum_{j \in [j^*+1,d]} \sigma_j^p - \sum_{j \in [k]} \sigma_j^p \left(\mathbf{W}^{\top} \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \right) \end{aligned}$$

Next, we show that for all $j \in [k]$, $\sigma_j(\mathbf{W}^\top(\mathbf{A} - \mathbf{A}_{j^*})) \ge \sigma_{j+j^*}(\mathbf{W}^\top \mathbf{A})$. Here, we invoke Fact 4.5 for $\mathbf{X} = (\mathbf{A} - \mathbf{A}_{j^*})$ and $\mathbf{Y} = \mathbf{A}_{j^*}$, with i = j and $j = j^*$. Note, the precondition on the indices i, j in Fact 4.5 is satisfied since \mathbf{X} , \mathbf{Y} are $n \times k$ matrices, and $j \in [k]$ and $j^* < k$. Then, we have

$$\sigma_{j+j^*} \left(\mathbf{W}^{\top} \mathbf{A} \right) = \sigma_{j+j^*} \left(\mathbf{W}^{\top} \left(\mathbf{A} - \mathbf{A}_{j^*} \right) + \mathbf{W}^{\top} \mathbf{A}_{j^*} \right)$$

$$\leq \sigma_j \left(\mathbf{W}^{\top} \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \right) + \sigma_{j^*+1} \left(\mathbf{W}^{\top} \mathbf{A}_{j^*} \right) ,$$

but $\mathbf{W}^{\top}\mathbf{A}_{j*}$ is a rank $\leq j^*$ matrix, and thus $\sigma_{j^*+1}\left(\mathbf{W}^{\top}\mathbf{A}_{j*}\right)=0$. Therefore, we can conclude,

$$\left\| \left(\mathbf{A} - \mathbf{A}_{j^*} \right) \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\mathsf{T}} \right) \right\|_{\mathcal{S}_p}^p \le \sum_{j \in [j^* + 1, d]} \sigma_j^p - \sum_{j \in [j^*, k + j^*]} \sigma_j^p \left(\mathbf{W}^{\mathsf{T}} \mathbf{A} \right)$$
(5.23)

Recall, for all $j \in [k]$, it follows from Equation (5.4) in the proof of Lemma 5.7 that $\sigma_j^p(\mathbf{AZ}) \ge \sigma_j^p(\mathbf{A}) - O(\gamma p) \sigma_{k+1}^2 \sigma_j^{p-2}$. Further, by definition, for $j \in [j^* + 1, k + j^*]$, $\sigma_j \le (1 + \epsilon/d) \sigma_{k+1}$ and thus, for all $j \in [j^* + 1, k]$,

$$\sigma_{j}^{p}(\mathbf{AZ}) \ge \sigma_{j}^{p} - O(\gamma p (1 + \epsilon/d)^{p-2}) \sigma_{k+1}^{p}$$

$$\ge \sigma_{j}^{p} - O(\gamma p) \sigma_{k+1}^{p}, \tag{5.24}$$

where the last inequality uses that $p = O(\log(d)/\epsilon)$. Finally, it follows from Lemma 5.9 that $\sigma_i^p(\mathbf{W}^{\mathsf{T}}\mathbf{A}) \geq \sigma_i^p(\mathbf{A}\mathbf{Z})$. Substituting this back into Equation (5.23), we have

$$\|\left(\mathbf{A} - \mathbf{A}_{j^*}\right) \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top}\right)\|_{\mathcal{S}_p}^p \leq \sum_{j \in [j^*+1,d]} \sigma_j^p - \sum_{j \in [j^*+1,k]} \sigma_j^p + O(\gamma k p) \sigma_{k+1}^p$$

$$\leq \left(1 + O(\gamma p k)\right) \|\mathbf{A} - \mathbf{A}_k\|_{\mathcal{S}_p}^p.$$
(5.25)

Taking the p-th root and substituting back into Equation (5.22),

$$\|\mathbf{A}\left(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top}\right)\|_{\mathcal{S}_{p}} \leq \left(1 + O(\gamma p k)\right)^{1/p} \|\mathbf{A} - \mathbf{A}_{k}\|_{\mathcal{S}_{p}} + O\left(\frac{\epsilon}{d}\right) \|\mathbf{A} - \mathbf{A}_{k}\|_{\mathcal{S}_{p}}, \tag{5.26}$$

and since $\gamma = \text{poly}(\epsilon/d)$, we have $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p} \le (1 + O(\epsilon)) \|\mathbf{A} - \mathbf{A}_k\|_{\mathcal{S}_p}$, which completes the analysis for this case.

Next, we consider the case where the gap between the top and the (k+1)-st singular value is small, i.e., $\sigma_1 < (1+1/p) \, \sigma_{k+1}$. We yet again split into cases, and consider the case where the Schatten-p norm of the tail is small, i.e. $\|\mathbf{A} - \mathbf{A}_k\|_{\mathcal{S}_p}^p \le \frac{k}{p^{1/3}\epsilon^{1/3}} \cdot \sigma_{k+1}^p$. Observe, for any $t \in [1, d-k-1]$,

$$\frac{k}{p^{1/3}\epsilon^{1/3}} \cdot \sigma_{k+1}^{p} \ge \|\mathbf{A} - \mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} \ge \sum_{i=k+1}^{k+1+t} \sigma_{i}^{p} \ge t \sigma_{k+1+t}^{p}.$$
 (5.27)

Then, setting $t = \frac{(1+1/p)^p k}{\epsilon^{1/3}p^{1/3}} = \Theta\left(\frac{k}{\epsilon^{1/3}p^{1/3}}\right)$, we have $\sigma_{k+1+t} \leq \sigma_{k+1}/(1+1/p)$. It suffices to show that we can detect this gap for some $s \geq k+1+t$. Recall, we know that $\hat{\sigma}_{k+1} = (1\pm 0.1/p)\sigma_{k+1}$ and $\hat{\sigma}_s = (1\pm 0.1/p)\sigma_s$. Then, we have

$$\hat{\sigma}_s \le \left(1 + \frac{0.1}{p}\right) \sigma_s \le \left(1 + \frac{0.1}{p}\right) \sigma_{k+1+t} \le \left(1 + \frac{0.1}{p}\right) \cdot \left(\frac{1}{1 + 1/p}\right) \sigma_{k+1} \le \frac{1}{\left(1 + \frac{0.5}{p}\right)} \hat{\sigma}_{k+1}. \tag{5.28}$$

Therefore, Algorithm 5.3 outputs **Z**, an orthonormal basis for $\mathbf{A}^{\mathsf{T}}\mathbf{W}_2$, where \mathbf{W}_2 is obtained by running Algorithm 5.6 on \mathbf{A}^{T} , initialized with a block size of $\Theta\left(\frac{k}{\epsilon^{1/3}p^{1/3}}\right)$ and run for $O(\log(d/\epsilon)\sqrt{p})$ iterations. Observe, since $\sigma_{k+1+t} \leq \sigma_{k+1}/(1+1/p)$, this suffices to demonstrate a gap that depends on p as follows: $\frac{\sigma_k}{\sigma_k - \sigma_{k+t+1}} \leq p$. Recall, we account for this gap by running $O(\log(d)\sqrt{p})$ iterations. Using the gap dependent analysis (Lemma 5.4), we can conclude that with probability at least 99/100, for all $i \in [k]$,

$$\left\|\mathbf{A}^{\top}(\mathbf{W}_{2})_{*,i}\right\|_{2}^{2} \ge \sigma_{i}^{2} - \operatorname{poly}\left(\frac{\epsilon}{d}\right)\sigma_{k+1}^{2}.$$
(5.29)

Then, applying Lemma 5.7 with $\mathbf{W}_2\mathbf{W}_2^{\mathsf{T}}$ satisfying the guarantee in (5.29), we have

$$\|\mathbf{A}^{\top}\mathbf{W}_{2}\mathbf{W}_{2}^{\top}\|_{\mathcal{S}_{p}}^{p} \geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - \operatorname{poly}\left(\frac{\epsilon}{d}\right) \sum_{i \in [k]} \sigma_{k+1}^{2} \sigma_{i}^{p-2}$$

$$\geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - \operatorname{poly}\left(\frac{\epsilon}{d}\right) \sigma_{k+1}^{p}.$$
(5.30)

where the last inequality uses that $\sigma_1 < (1+1/p)\sigma_{k+1}$ and $(1+1/p)^{p-2} = O(1)$. Next, we use Lemma 5.5 to relate $\|\mathbf{A}^{\top}\mathbf{W}_2\mathbf{W}_2^{\top}\|_{\mathcal{S}_p}^p$ to $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p}^p$, where \mathbf{Z} is an orthonormal basis for $\mathbf{A}^{\top}\mathbf{W}_2\mathbf{W}_2^{\top}$ as output by the algorithm. Setting $\mathbf{Q} = \mathbf{Z}\mathbf{Z}^{\top}$ and $\mathbf{P} = \mathbf{W}_2\mathbf{W}_2^{\top}$, we observe that $\|\mathbf{P}\mathbf{A}\mathbf{Q}\|_{\mathcal{S}_p}^p = \|\mathbf{A}^{\top}\mathbf{W}_2\mathbf{W}_2^{\top}\|_{\mathcal{S}_p}^p = \|\mathbf{W}_2\mathbf{W}_2^{\top}\mathbf{A}\|_{\mathcal{S}_p}^p$ and $\|(\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_p}^p = \|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_p}^p$. Then, invoking Lemma 5.5 and plugging in Equation (5.30), we have

$$\|(\mathbf{I} - \mathbf{P}) \mathbf{A} (\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{A} (\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_{p}}^{p} \leq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}^{\top}\mathbf{W}_{2}\mathbf{W}_{2}^{\top}\|_{\mathcal{S}_{p}}^{p}$$

$$\leq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} + \operatorname{poly}\left(\frac{\epsilon}{d}\right)\sigma_{k+1}^{p}$$

$$\leq \left(1 + \operatorname{poly}\left(\frac{\epsilon}{d}\right)\right)\|\mathbf{A} - \mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p},$$
(5.31)

which concludes the analysis in this case.

As shown in Equation 5.28, we can detect a gap between σ_{k+1+t} and σ_{k+1} by comparing $\hat{\sigma}_s$ and $\hat{\sigma}_{k+1}$. When 5.28 does not hold, we know that $\hat{\sigma}_s \geq (1 + 0.5/p) \hat{\sigma}_{k+1}$ and Algorithm 5.3 outputs **Z**, an orthonormal basis for $\mathbf{A}^{\mathsf{T}}\mathbf{W}_1\mathbf{W}_1^{\mathsf{T}}$. Since we have $(1 \pm 0.1/p)$ -approximate estimates to these quantities, we can conclude that $\sigma_s \geq (1 + 0.1/p) \sigma_{k+1}$. Then, we have

$$\|\mathbf{A} - \mathbf{A}_k\|_{\mathcal{S}_p}^p \ge s \cdot \sigma_s^p = \Omega\left(\frac{k}{\epsilon^{1/3}p^{1/3}}\right)\sigma_{k+1}^p.$$

It therefore remains to consider the case where $\|\mathbf{A} - \mathbf{A}_k\|_{\mathcal{S}_p}^p > \frac{ck}{p^{1/3}\epsilon^{1/3}} \cdot \sigma_{k+1}^p$, for a fixed universal constant c. Here, we note that the tail is large enough that an additive error of $O(\epsilon^{2/3}p^{1/3})$ σ_{k+1}^2 on each of the top-k singular values suffices. Formally, it follows from Lemma 5.2 (setting $\gamma = \epsilon^{2/3}p^{-1/3}$, and invoking it for \mathbf{A}^{\top}) that initializing Algorithm 5.6 with block size k and running for $O(\log(d/\epsilon)p^{1/6}/\epsilon^{1/3})$ iterations suffices to output a $n \times k$ matrix \mathbf{W}_1 such that with probability at least 99/100, for all $i \in [k]$,

$$\|\mathbf{A}^{\top}(\mathbf{W}_{1})_{*,i}\|_{2}^{2} \geq \sigma_{i}^{2} - \epsilon^{2/3} p^{-1/3} \sigma_{k+1}^{2}.$$

Then, invoking Lemma 5.7 with \mathbf{A}^{T} and \mathbf{W}_1 as defined above, we have

$$\|\mathbf{A}^{\top}\mathbf{W}_{1}\mathbf{W}_{1}^{\top}\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{W}_{1}\mathbf{W}_{1}^{\top}\mathbf{A}\|_{\mathcal{S}_{p}}^{p}$$

$$\geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - \sum_{i \in [k]} O\left(\epsilon^{2/3}p^{-1/3}p\right)\sigma_{k+1}^{2}\sigma_{i}^{p-2}$$

$$\geq \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} - O\left(k\epsilon^{2/3}p^{2/3}\right)\sigma_{k+1}^{p}$$

$$(5.32)$$

where the last inequality uses that $\sigma_1 < (1 + 1/p)\sigma_{k+1}$ and $(1 + 1/p)^p = O(1)$. Recall, in this case, Algorithm 5.3 outputs $\mathbf{Z}\mathbf{Z}^{\mathsf{T}}$ where \mathbf{Z} is an orthonormal basis for $\mathbf{A}^{\mathsf{T}}\mathbf{W}_1\mathbf{W}_1^{\mathsf{T}}$. Next, we invoke Lemma

5.5 to relate $\|\mathbf{A}^{\mathsf{T}}\mathbf{W}_{1}\mathbf{W}_{1}^{\mathsf{T}}\|_{\mathcal{S}_{p}}^{p}$ to $\|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\mathsf{T}})\|_{\mathcal{S}_{p}}^{p}$. Setting $\mathbf{Q} = \mathbf{Z}\mathbf{Z}^{\mathsf{T}}$ and $\mathbf{P} = \mathbf{W}_{1}\mathbf{W}_{1}^{\mathsf{T}}$, we observe that $\|\mathbf{P}\mathbf{A}\mathbf{Q}\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{W}_{1}\mathbf{W}_{1}^{\mathsf{T}}\mathbf{A}\|_{\mathcal{S}_{p}}^{p}$ and $\|(\mathbf{I} - \mathbf{P})\mathbf{A}(\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{A}(\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\mathsf{T}})\|_{\mathcal{S}_{p}}^{p}$. Then, invoking Lemma 5.5 and plugging in Equation (5.32), we have

$$\|(\mathbf{I} - \mathbf{P}) \mathbf{A} (\mathbf{I} - \mathbf{Q})\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{A} (\mathbf{I} - \mathbf{Z}\mathbf{Z}^{\top})\|_{\mathcal{S}_{p}}^{p} \leq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{W}_{1}\mathbf{W}_{1}^{\top}\mathbf{A}\|_{\mathcal{S}_{p}}^{p}$$

$$\leq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p} + O(k\epsilon^{2/3}p^{2/3})\sigma_{k+1}^{p}$$

$$\leq (1 + O(p\epsilon)) \|\mathbf{A} - \mathbf{A}_{k}\|_{\mathcal{S}_{p}}^{p},$$
(5.33)

where the last inequality follows from our assumption on the Schatten-p norm of the tail, given the case we are in. Taking the (1/p)-th root, and recalling that $\epsilon < 1/2$, we obtain

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\mathsf{T}} \right) \right\|_{\mathcal{S}_{p}} \le (1 + O(\epsilon)) \left\| \mathbf{A} - \mathbf{A}_{k} \right\|_{p}, \tag{5.34}$$

which concludes the final case.

Next, we analyze the running time and matrix-vector products. Running Algorithm 5.6 with block size k for $q = O(\log(d)p^{1/6}/\epsilon^{1/3})$ iterations requires $O\left(\frac{\operatorname{nnz}(\mathbf{A})kp^{1/6}\log(d)}{\epsilon^{1/3}}\right)$ time and $O\left(\frac{kp^{1/6}\log(d)}{\epsilon^{1/3}}\right)$ matrix-vector products. Similarly, running with block size $O\left(k/(\epsilon p)^{1/3}\right)$ for $q = O(\log(d/\epsilon)\sqrt{p})$ iterations requires $O\left(\frac{\operatorname{nnz}(\mathbf{A})kp^{1/6}\log(d/\epsilon)}{\epsilon^{1/3}}\right)$ time and $O\left(\frac{kp^{1/6}\log(d)}{\epsilon^{1/3}}\right)$ matrix-vector products. Finally, we observe that to obtain a (1+1/p)-approximation to σ_1 and σ_{k+1} , we need $O(\log(d)\sqrt{p})$ iterations with blocksize k+1 and this requires $O(\log(d)\sqrt{p}k)$ matrix-vector products. Note, our setting of the exponent of p and ϵ was chosen to balance the two cases, and this concludes the proof.

6 Query Lower Bounds

Next, we show that the ϵ -dependence obtained by our algorithms for Schatten-p low-rank approximation is optimal in the restricted computation model of matrix-vector products. The matrix-vector product model is defined as follows: given a matrix \mathbf{A} , our algorithm is allowed to make adaptive matrix-vector queries to \mathbf{A} , where one matrix-vector query is of the form $\mathbf{A}v$, for any $v \in \mathbb{R}^d$. Our lower bounds are information-theoretic and rely on the hardness of estimating the smallest eigenvalue of a Wishart ensemble, as established in recent work of Braverman, Hazan, Simchowitz and Woodworth [BHSW20].

We split the lower bounds into the case of $p \in [1,2]$ and p > 2. For $p \in [1,2]$, we have a simple argument based on the Araki-Lieb-Thirring inequality (Fact 4.10), whereas for p > 2, our lower bounds require an involved argument using a norm compression inequality for partitioned operators (Fact 4.14).

6.1 Lower Bounds for $p \in [1, 2]$

The main lower bound we prove in this sub-section is as follows:

Theorem 6.1 (Query Lower Bound for $p \in [1,2]$). Given $\varepsilon > 0$, and $p \in [1,2]$, there exists a distribution \mathcal{D} over $n \times n$ matrices such that for $\mathbf{A} \sim \mathcal{D}$, any randomized algorithm that with probability at least 9/10 outputs a rank-1 matrix \mathbf{B} such that $\|\mathbf{A} - \mathbf{B}\|_{\mathcal{S}_p}^p \leq (1 + \varepsilon) \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p$ must make $\Omega(1/\varepsilon^{1/3})$ matrix-vector queries to \mathbf{A} .

We require the following theorem on the hardness of computing the minimum eigenvalue of a Wishart Matrix, introduced recently by Braverman, Hazan, Simchowitz and Woodworth [BHSW20]:

Theorem 6.2 (Computing Min Eigenvalue of Wishart, Theorem 3.1 [BHSW20]). *Given* $\epsilon \in (0,1)$, there exists a function $\mathbf{d} : (0,1) \to \mathbb{N}$ such that for all $d \geq \mathbf{d}(\epsilon)$, the following holds. Let $\mathbf{W} \sim \text{Wishart}(d)$ be a Wishart matrix and $\{\lambda_i\}_{i\in[d]}$ be the eigenvalues of \mathbf{W} , in descending order. Then, there exists a universal constant c^* such that:

- 1. Let ζ_1 be the event that $\lambda_d(\mathbf{W}) \leq c_1/d^2$, ζ_2 be the event that $\lambda_{d-1}(\mathbf{W}) \lambda_d(\mathbf{W}) \geq c_2/d^2$ and ζ_3 be the event that $\|\mathbf{W}\|_{op} \leq 5$, where c_1 and c_2 are constants that depend only on ϵ . Then, $\Pr_{\mathbf{W}}[\zeta_1 \cap \zeta_2 \cap \zeta_3] \geq 1 \frac{c^* \sqrt{\epsilon}}{2}$.
- 2. Any randomized algorithm that makes at most $(1 \epsilon)d$ adaptive matrix-vector queries and outputs an estimate $\hat{\lambda}_d$ must satisfy

$$\Pr_{\mathbf{W}} \left| \left| \hat{\lambda}_d - \lambda_d \right| \ge \frac{1}{4d^2} \right| \ge c^* \sqrt{\epsilon}.$$

We also use the following lemma from [BHSW20] bounding the minimum eigenvalue of a Wishart ensemble:

Lemma 6.3 (Non-Asymptotic Spectra of Wishart Ensembles, Corollary 3.3 [BHSW20]). Let $\mathbf{W} \sim \text{Wishart}(n)$ be such that $n = \Omega(1/\varepsilon^3)$. Then, there exists a universal constant $c_2 > 0$ such that

$$\Pr\left[\lambda_n\left(\mathbf{W}\right) \ge \frac{1}{n^2}\right] \ge c_2, \quad and \quad \Pr\left[\lambda_n\left(\mathbf{W}\right) < \frac{1}{2n^2}\right] \ge \frac{c_2}{2}.$$

We are now ready to prove Theorem 6.1. Our high level approach is to show that we can take any solution that is a $(1 + \varepsilon)$ -relative-error Schatten-p low-rank approximation to the hard instance $I - \frac{1}{5}W$, where W is a Wishart ensemble, and extract from it an accurate estimate of the minimum eigenvalue of W, thus appealing to the hardness stated in (2) of Theorem 6.2 above.

Proof of Theorem 6.1. Let $n = \Theta\left(1/\epsilon^{1/3}\right)$ and let $\mathbf{A} = \mathbf{I} - \frac{1}{5}\mathbf{W}$ be an $n \times n$ instance where $\mathbf{W} \sim \mathbf{W}$ is with ζ_1 be the event that $\|\mathbf{W}\|_{\mathrm{op}} \leq 5$. It follows from Fact 4.16 that ζ_1 holds with probability at least 99/100, and we condition on this event. Let ζ_2 be the event that $\lambda_n(\mathbf{W}) \geq \frac{1}{n^2} = \frac{\epsilon^{2/3}}{\epsilon^*}$ and ζ_3 be the event that $\lambda_n(\mathbf{W}) < \frac{1}{2n^2} = \frac{\epsilon^{2/3}}{2c^*}$.

Then, conditioning on ζ_2 , we have that

$$1 - \frac{1}{5}\lambda_n(\mathbf{W}) \le 1 - \frac{\varepsilon^{2/3}}{5c^*}.\tag{6.1}$$

Similarly, conditioning on ζ_3 , we have that

$$1 - \frac{1}{5}\lambda_n(\mathbf{W}) \ge 1 - \frac{\varepsilon^{2/3}}{10c^*}.$$
 (6.2)

We observe that for $p \in [1, 2]$, using Bernoulli's inequality (Fact 4.6) we have

$$\left(1 - \frac{1}{5}\lambda_n(\mathbf{W})\right)^p \ge 1 - \frac{p}{5}\lambda_n(\mathbf{W})$$

and since $(1-x)^p \le (1-x)$ for any $x \in (0,1)$, we also have that,

$$\left(1 - \frac{1}{5}\lambda_n(\mathbf{W})\right)^p \le 1 - \frac{1}{5}\lambda_n(\mathbf{W})$$

Therefore, we can conclude, $\left(1 - \frac{1}{5}\lambda_n(\mathbf{W})\right)^p = 1 - \Theta\left(\lambda_n(\mathbf{W})\right)$. Further, it follows from part (1) of Fact 4.16 that $0 \le \mathbf{I} - \frac{1}{5}\mathbf{W} \le \mathbf{I}$, and thus

$$\|\mathbf{A}\|_{\mathcal{S}_p}^p = \sum_{i \in [n]} \lambda_i^p \left(\mathbf{I} - \frac{1}{5} \mathbf{W} \right) \le \sum_{i \in [n]} \lambda_i \left(\mathbf{I} - \frac{1}{5} \mathbf{W} \right) \le O\left(\frac{1}{\epsilon^{1/3}}\right)$$
(6.3)

where the last inequality follows from the fact that $n = \sqrt{c^*}/\epsilon^{1/3}$. Let \mathbf{A}_1 denote the best rank-1 approximation to \mathbf{A} . Then, it follows from Equation (6.3) that

$$\epsilon \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p \le \epsilon \|\mathbf{A}\|_{\mathcal{S}_p}^p \le O(\epsilon^{2/3})$$
 (6.4)

Observe, any $(1 + \epsilon)$ -approximate relative-error Schatten-p low-rank approximation algorithm for k = 1 outputs a matrix vv^{T} such that

$$\|\mathbf{A} \left(\mathbf{I} - vv^{\top}\right)\|_{\mathcal{S}_{p}}^{p} \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A}_{1}\|_{\mathcal{S}_{p}}^{p}$$

$$\leq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}\|_{\mathrm{op}}^{p} + \Theta(\epsilon^{2/3})$$
(6.5)

By definition of the Schatten-*p* norm we have:

$$\|\mathbf{A} \left(\mathbf{I} - vv^{\top}\right)\|_{\mathcal{S}_{p}}^{p} = \operatorname{Tr}\left(\left(\left(\mathbf{I} - vv^{\top}\right)^{2} \mathbf{A}^{2} \left(\mathbf{I} - vv^{\top}\right)^{2}\right)^{p/2}\right)$$

$$\geq \operatorname{Tr}\left(\left(\mathbf{I} - vv^{\top}\right)^{p} \mathbf{A}^{p} \left(\mathbf{I} - vv^{\top}\right)^{p}\right)$$

$$= \operatorname{Tr}\left(\mathbf{A}^{p} - \mathbf{A}^{p}vv^{\top}\right)$$

$$= \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \operatorname{Tr}\left(\left(vv^{\top}\right)^{p/2} \left(\mathbf{A}^{2}\right)^{p/2} \left(vv^{\top}\right)^{p/2}\right)$$

$$\geq \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \operatorname{Tr}\left(\left(vv^{\top}\mathbf{A}^{2}vv^{\top}\right)^{p/2}\right)$$

$$= \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}vv^{\top}\|_{\mathcal{S}_{p}}^{p}$$

$$= \|\mathbf{A}\|_{\mathcal{S}_{p}}^{p} - \|\mathbf{A}v\|_{2}^{p}$$

$$(6.6)$$

where the first and last inequality follows from the reverse Araki-Lieb-Thirring inequality (Fact 4.10). Combining equations (6.5) and (6.6), we have that

$$\|\mathbf{A}\|_{\text{op}}^{p} \ge \|\mathbf{A}v\|_{2}^{p} \ge \|\mathbf{A}\|_{\text{op}}^{p} - \Theta(\epsilon^{2/3})$$
 (6.7)

Next, we observe that $\mathbf{A}v = (\mathbf{I} - 1/5\mathbf{W}) v$ can be computed with one additional matrix-vector product and

$$\|\mathbf{A}\|_{\text{op}}^{p} = \left(1 - \frac{1}{5}\lambda_{n}(\mathbf{W})\right)^{p} = 1 - \frac{p}{5}\lambda_{n}(\mathbf{W}) + O(\lambda_{n}^{2}(\mathbf{W}))$$
 (6.8)

Consider the estimator $\hat{\lambda}(\mathbf{W}) = \frac{5}{p} \left(1 - \left\| \left(\mathbf{I} - \frac{1}{5} \mathbf{W} \right) v \right\|_{2}^{p} \right)$. Combining equations (6.7) and (6.8), we can conclude

$$\hat{\lambda}(\mathbf{W}) = \lambda_{\min}(\mathbf{W}) \pm \Theta(\epsilon^{2/3}).$$

obtaining an additive error estimate to the minimum eigenvalue of **W** by computing an additional matrix-vector product. It follows that we satisfy conditions (1) and (2) in Theorem 6.2 and thus any algorithm for computing a rank-1 approximation to the matrix $\mathbf{A} = \mathbf{I} - \frac{1}{5}\mathbf{W}$ in Schatten p norm must make at least $\frac{1}{\varepsilon^{1/3}}$ queries to the aforementioned matrix, completing the proof. The claim follows from Theorem 6.2.

6.2 Lower Bound for p > 2

We now consider the case when p > 2. We note that the previous approach no longer works since we cannot lower bound the cost of $\| (\mathbf{I} - \mathbf{W}/5) (\mathbf{I} - vv^{\mathsf{T}}) \|_{\mathcal{S}_p}$, as the Araki-Lieb-Thirring inequality reverses (see application in Equation 6.6). Therefore, we require a new approach, and appeal to a special case of Conjecture 2.2 that is known to be true, i.e. the Aligned Norm Compression inequality (see Fact 4.14). The main theorem we prove in this sub-section is as follows:

Theorem 6.4 (Query Lower Bound for p > 2). Given $\varepsilon > 0$, and $p \ge 2$ such that p = O(1), there exists a distribution \mathcal{D} over $n \times n$ matrices such that for $\mathbf{A} \sim \mathcal{D}$, any randomized algorithm that with probability at least 99/100 outputs a unit vector u such that $\|\mathbf{A} - \mathbf{A}uu^{\top}\|_{\mathcal{S}_p}^p \le (1 + \varepsilon) \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p$ must make $\Omega\left(1/\varepsilon^{1/3}\right)$ matrix-vector queries to \mathbf{A} .

We first introduce a sequence of key lemmas required for our proof.

Corollary 6.5 (Special Case of Lemma 5.2). Given $\gamma \in [0,1]$, a vector $v \in \mathbb{R}^d$ and an $n \times d$ matrix \mathbf{A} , let $t = \log(n/\gamma)/(c\sqrt{\gamma})$, for a fixed universal constant c. Then, there exists an algorithm that computes t matrix-vector products with \mathbf{A} and outputs a unit vector u such that with probability at least 99/100,

$$\|\mathbf{A}\|_{op}^2 - \|\mathbf{A}u\|_2^2 \le O(\gamma \sigma_2^2).$$

where σ_2 is the second largest singular value of **A**.

Next, we prove a key lemma relating the norm of a matrix to norms of orthogonal projections applied to the matrix. We note that this lemma is straight forward and holds for arbitrary vectors unit u, v if Conjecture 2.2 holds. However, we show that we can transform our matrix to have structure such that we can apply Fact 4.14 instead.

Lemma 6.6 (Orthogonal Projectors to Block Matrices). Given an $n \times d$ matrix \mathbf{A} , p > 2 and unit vectors $u \in \mathbb{R}^d$, $v \in \mathbb{R}^n$, such that $(\mathbf{I} - vv^\top) \mathbf{A} u u^\top = 0$. Then, we have

$$\|\mathbf{A}\|_{\mathcal{S}_p} \leq \left\| \begin{pmatrix} \|vv^{\top}\mathbf{A}uu^{\top}\|_{\mathcal{S}_p} & \|vv^{\top}\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p} \\ 0 & \|(\mathbf{I} - vv^{\top})\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p} \end{pmatrix} \right\|_{\mathcal{S}_n}.$$

Proof. Let $\mathbf{I} - vv^{\top} = \mathbf{Y}\mathbf{Y}^{\top}$, where \mathbf{Y} has n-1 orthonormal columns. Further, since v and \mathbf{Y} span disjoint subspaces, and the union of their span is \mathbb{R}^n , the matrix $(v \mid \mathbf{Y})$, obtained by concatenating their columns is unitary. Then, let $\mathbf{R} = (v \mid \mathbf{Y})^{\top}$ and observe, \mathbf{R} has orthonormal rows and columns (since \mathbf{R} is unitary). Next, let $\mathbf{I} - uu^{\top} = \mathbf{Z}\mathbf{Z}^{\top}$, where \mathbf{Z} is $d \times (d-1)$ and has orthonormal columns. Let $\mathbf{S} = (u \mid \mathbf{Z})^{\top}$, and observe \mathbf{S} has orthonormal rows and columns.

Let $\hat{\mathbf{A}} = \mathbf{R}\mathbf{A}\mathbf{S}^{\mathsf{T}}$, which admits the following block-matrix form:

$$\hat{\mathbf{A}} = \begin{pmatrix} v^{\top} \\ \mathbf{Y}^{\top} \end{pmatrix} \cdot \mathbf{A} \cdot (u \mid \mathbf{Z}) = \begin{pmatrix} v^{\top} \\ \mathbf{Y}^{\top} \end{pmatrix} (\mathbf{A}u \mid \mathbf{A}\mathbf{Z}) = \begin{pmatrix} v^{\top} \mathbf{A}u & v^{\top} \mathbf{A}\mathbf{Z} \\ \mathbf{Y}^{\top} \mathbf{A}u & \mathbf{Y}^{\top} \mathbf{A}\mathbf{Z} \end{pmatrix}$$

Since **R** and **S** are unitary, it follows from unitary invariance of the Schatten-*p* norm that

$$\|\mathbf{A}\|_{\mathcal{S}_p} = \|\hat{\mathbf{A}}\|_{\mathcal{S}_p} = \left\| \begin{pmatrix} v^{\top} \mathbf{A} u & v^{\top} \mathbf{A} \mathbf{Z} \\ \mathbf{Y}^{\top} \mathbf{A} u & \mathbf{Y}^{\top} \mathbf{A} \mathbf{Z} \end{pmatrix} \right\|_{\mathcal{S}_p} = \left\| \begin{pmatrix} v^{\top} \mathbf{A} u & v^{\top} \mathbf{A} \mathbf{Z} \\ 0 & \mathbf{Y}^{\top} \mathbf{A} \mathbf{Z} \end{pmatrix} \right\|_{\mathcal{S}_p}, \tag{6.9}$$

where the last equality follows from observing that $\|\mathbf{Y}^{\mathsf{T}}\mathbf{A}u\|_F = \|\mathbf{Y}\mathbf{Y}^{\mathsf{T}}\mathbf{A}uu^{\mathsf{T}}\|_F = \|(\mathbf{I} - vv^{\mathsf{T}})\mathbf{A}uu^{\mathsf{T}}\|_F = 0$ and therefore $\mathbf{Y}^{\mathsf{T}}\mathbf{A}u$ is a matrix of all 0s. Next, we append a set of d-2 columns of 0's to make the top left and top right block the same size. Since this does not change the singular values, we have

$$\|\mathbf{A}\|_{\mathcal{S}_p} = \left\| \begin{pmatrix} v^{\mathsf{T}} \mathbf{A} u & 0 & v^{\mathsf{T}} \mathbf{A} \mathbf{Z} \\ 0 & 0 & \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \end{pmatrix} \right\|_{\mathcal{S}_p}$$
(6.10)

Next, we construct a rotation matrix \mathbf{R} such that on right multiplying a row vector by \mathbf{R} , the first d-1 coordinates remain the same and on the remaining coordinates, the vector $v^{\mathsf{T}}\mathbf{AZ}$ gets mapped to ce_1^{T} for some scalar c. Let \mathbf{S} be the $d-1\times d-1$ rotation matrix such that $v^{\mathsf{T}}\mathbf{AZS}=ce_1^{\mathsf{T}}$. Then, $\mathbf{R}=\begin{pmatrix}\mathbf{I}&0\\0&\mathbf{S}\end{pmatrix}$ and it is easy to verify that \mathbf{R} is unitary. Therefore,

$$\begin{pmatrix} v^{\mathsf{T}} \mathbf{A} u & 0 & v^{\mathsf{T}} \mathbf{A} \mathbf{Z} \\ 0 & 0 & \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \end{pmatrix} \cdot \mathbf{R} = \begin{pmatrix} v^{\mathsf{T}} \mathbf{A} u & 0 & c e_1^{\mathsf{T}} \\ 0 & 0 & \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \mathbf{S} \end{pmatrix}$$

Now, we observe the final matrix above has a block matrix form we can apply the Aligned Norm Compression inequality from Fact 4.14, with $\alpha_1 = v^T \mathbf{A}u$, $\alpha_2 = c$, $\beta_1 = 0$ and $\beta_2 = 0$, and therefore

$$\|\mathbf{A}\|_{\mathcal{S}_{p}} = \left\| \begin{pmatrix} v^{\mathsf{T}} \mathbf{A} u & 0 & c e_{1}^{\mathsf{T}} \\ 0 & 0 & \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \mathbf{S} \end{pmatrix} \right\|_{\mathcal{S}_{p}} \leq \left\| \begin{pmatrix} \|v^{\mathsf{T}} \mathbf{A} u\|_{\mathcal{S}_{p}} & 0 & \|c e_{1}^{\mathsf{T}}\|_{\mathcal{S}_{p}} \\ 0 & 0 & \|\mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \mathbf{S}\|_{\mathcal{S}_{p}} \end{pmatrix} \right\|_{\mathcal{S}_{p}}$$

$$= \left\| \begin{pmatrix} \|vv^{\mathsf{T}} \mathbf{A} uu^{\mathsf{T}}\|_{\mathcal{S}_{p}} & \|vv^{\mathsf{T}} \mathbf{A} \mathbf{Z} \mathbf{Z}^{\mathsf{T}}\|_{\mathcal{S}_{p}} \\ 0 & \|\mathbf{Y} \mathbf{Y}^{\mathsf{T}} \mathbf{A} \mathbf{Z} \mathbf{Z}^{\mathsf{T}}\|_{\mathcal{S}_{p}} \end{pmatrix} \right\|_{\mathcal{S}_{p}}$$

$$(6.11)$$

where the last equality follows from unitary invariance and substituting the definition of \mathbf{YY}^{T} and \mathbf{ZZ}^{T} completes the proof.

Fact 6.7 (SVD of a 2 × 2 Matrix). Given a 2 × 2 matrix $\mathbf{M} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ let $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\mathsf{T}}$ be the SVD of \mathbf{M} . Then,

$$\Sigma_{1,1} = \sqrt{\frac{a^2 + b^2 + c^2 + d^2 + \sqrt{(a^2 + b^2 - c^2 - d^2)^2 + 4(ac + bd)^2}}{2}},$$

and

$$\Sigma_{2,2} = \sqrt{\frac{a^2 + b^2 + c^2 + d^2 - \sqrt{(a^2 + b^2 - c^2 - d^2)^2 + 4(ac + bd)^2}}{2}}.$$

Now, we are ready to prove Theorem 6.4.

Proof of Theorem 6.4. Let $\mathbf{A} = \mathbf{I} - \frac{1}{5}\mathbf{W}$ where \mathbf{W} is an $n \times n$ Wishart matrix as in the proof of Theorem 6.1 and we have by hypothesis that there is an algorithm that with probability at least 99/100, outputs a unit vector u such that $\|\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p}^p \leq (1 + \varepsilon) \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^p$. Let $v = \mathbf{A}u/\|\mathbf{A}u\|_2$ and observe, $(\mathbf{I} - vv^{\top})\mathbf{A}uu^{\top} = 0$. Further, by the unitary invariance of the Schatten-p norm,

$$\left\| vv^{\mathsf{T}} \mathbf{A} u u^{\mathsf{T}} \right\|_{\mathcal{S}_p} = \left| v^{\mathsf{T}} \mathbf{A} u \right| = \frac{\left| u^{\mathsf{T}} \mathbf{A}^{\mathsf{T}} \mathbf{A} u \right|}{\left\| \mathbf{A} u \right\|_2} = \left\| \mathbf{A} u \right\|_2. \tag{6.12}$$

Similarly,

$$\|vv^{\top}\mathbf{A}\left(\mathbf{I} - uu^{\top}\right)\|_{\mathcal{S}_{p}} = \sqrt{\|v^{\top}\mathbf{A}\left(\mathbf{I} - uu^{\top}\right)\|_{2}^{2}} = \sqrt{\|v^{\top}\mathbf{A}\|_{2}^{2} - \|v^{\top}\mathbf{A}uu^{\top}\|_{2}^{2}}$$

$$= \sqrt{\frac{\|u^{\top}\mathbf{A}^{\top}\mathbf{A}\|_{2}^{2}}{\|\mathbf{A}u\|_{2}^{2}} - \|\mathbf{A}u\|_{2}^{2}}$$

$$\leq \sqrt{\frac{\|u^{\top}\mathbf{A}^{\top}\|_{2}^{2} \cdot \|\mathbf{A}\|_{op}^{2}}{\|\mathbf{A}u\|_{2}^{2}} - \|\mathbf{A}u\|_{2}^{2}}$$

$$\leq \epsilon^{1/3}\sigma_{2},$$
(6.13)

where we use sub-multiplicativity of the ℓ_2 norm and Corollary 6.5 with $\gamma = \epsilon^{2/3}$. Note that we can assume w.l.o.g. that Corollary 6.5 holds since we can just iterate Block Krylov $q = (1/c\epsilon^{1/3})$ times, for a sufficiently large constant c, starting the iterations with the vector u output by the algorithm hypothesized for the theorem, and pay only $(1/c\epsilon^{1/3})$ extra matrix-vector products. Since $vv^{\mathsf{T}}\mathbf{A} + \mathbf{A}uu^{\mathsf{T}} - vv^{\mathsf{T}}\mathbf{A}uu^{\mathsf{T}}$ has rank at most 3,

$$\|(\mathbf{I} - vv^{\top}) \mathbf{A} (\mathbf{I} - uu^{\top})\|_{\mathcal{S}_{p}}^{p} = \|\mathbf{A} - vv^{\top}\mathbf{A} - \mathbf{A}uu^{\top} + vv^{\top}\mathbf{A}uu^{\top}\|_{\mathcal{S}_{p}}^{p}$$

$$\geq \|\mathbf{A} - \mathbf{A}_{3}\|_{\mathcal{S}_{p}}^{p}$$

$$= \Omega \left(\frac{1}{\epsilon^{1/3}}\right),$$
(6.14)

where the last inequality follows from Fact 4.16.

Let
$$\mathbf{M} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} = \begin{pmatrix} \|vv^{\top}\mathbf{A}uu^{\top}\|_{\mathcal{S}_p} & \|vv^{\top}\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p} \\ \|(\mathbf{I} - vv^{\top})\mathbf{A}uu^{\top}\|_{\mathcal{S}_p} & \|(\mathbf{I} - vv^{\top})\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p} \end{pmatrix}^{\top}$$
. Then, it follows from Fact 6.7 that

$$\Sigma_{1,1}(\mathbf{M}) = \frac{1}{\sqrt{2}} \cdot \sqrt{a^2 + c^2 + d^2 + \sqrt{(a^2 - c^2 - d^2)^2 + 4(ac)^2}}$$

$$= \frac{1}{\sqrt{2}} \cdot \sqrt{a^2 + c^2 + d^2 + (c^2 + d^2 - a^2) + \Theta\left(\frac{4a^2c^2}{c^2 + d^2 - a^2}\right)}$$

$$= \sqrt{c^2 + d^2 + \Theta\left(\frac{a^2c^2}{c^2 + d^2 - a^2}\right)},$$
(6.15)

where we use that b=0, c, $a \le 1$ and $1 \ll d$ and the Taylor expansion of $\sqrt{x+y}$ for x, $y \ge 0$. Similarly,

$$\Sigma_{2,2}(\mathbf{M}) = \sqrt{a^2 - \Theta\left(\frac{a^2c^2}{c^2 + d^2 - a^2}\right)}.$$
 (6.16)

Then, using equations (6.15) and (6.16) we can bound the Schatten-p norm of **M** as follows:

$$\|\mathbf{M}\|_{\mathcal{S}_{p}}^{p} \leq \underbrace{\left(c^{2} + d^{2} + \Theta\left(\frac{a^{2}c^{2}}{c^{2} + d^{2} - a^{2}}\right)\right)^{p/2}}_{6.17.1} + \underbrace{\left(a^{2} - \Theta\left(\frac{a^{2}c^{2}}{c^{2} + d^{2} - a^{2}}\right)\right)^{p/2}}_{6.17.2}.$$
(6.17)

We now bound each of the terms above. Consider the first term:

$$\left(c^{2} + d^{2} + \Theta\left(\frac{a^{2}c^{2}}{c^{2} + d^{2} - a^{2}}\right)\right)^{p/2} = \left(\left\|vv^{\mathsf{T}}\mathbf{A}\left(\mathbf{I} - uu^{\mathsf{T}}\right)\right\|_{\mathcal{S}_{p}}^{2} + \left\|\left(\mathbf{I} - vv^{\mathsf{T}}\right)\mathbf{A}\left(\mathbf{I} - uu^{\mathsf{T}}\right)\right\|_{\mathcal{S}_{p}}^{2} + \Theta\left(\varepsilon^{2/3}\left\|\mathbf{A}u\right\|_{2}^{2}\right)\right)^{p/2} \\
\leq \left(\Theta\left(\varepsilon^{2/3}\right) + \left\|\mathbf{A}\left(\mathbf{I} - uu^{\mathsf{T}}\right)\right\|_{\mathcal{S}_{p}}^{2}\right)^{p/2} \\
\leq \left(1 + O\left(\varepsilon^{2p/3}\right)\right) \left\|\mathbf{A} - \mathbf{A}_{1}\right\|_{\mathcal{S}_{p}}^{p}, \tag{6.18}$$

where we use equation (6.12), (6.13), and (6.14), and $\|\mathbf{A}(\mathbf{I} - uu^{\top})\|_{\mathcal{S}_p}^2 \le (1 + \epsilon)^{2/p} \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^2$. The last inequality follows from observing that

$$\varepsilon^{2/3} \leq O\left(\varepsilon^{4/3} \cdot \frac{1}{\varepsilon^{2/3p}}\right) \leq O\left(\varepsilon^{4/3} \cdot \|\mathbf{A} - \mathbf{A}_1\|_{\mathcal{S}_p}^2\right).$$

We can now bound the second term in Equation 6.17 as follows:

$$\left(a^{2} - \Theta\left(\frac{a^{2}c^{2}}{c^{2} + d^{2} - a^{2}}\right)\right)^{p/2} = \left(\|\mathbf{A}u\|_{2}^{2} - \Theta\left(\epsilon^{2/3}\|\mathbf{A}u\|_{2}^{2}\right)\right)^{p/2} \le \|\mathbf{A}u\|_{2}^{p}. \tag{6.19}$$

Then, we have

$$\|\mathbf{M}\|_{\mathcal{S}_{p}}^{p} \leq \left(1 + O\left(\epsilon^{2p/3}\right)\right) \|\mathbf{A} - \mathbf{A}_{1}\|_{\mathcal{S}_{p}}^{p} + \|\mathbf{A}u\|_{2}^{p}.$$

It follows from Lemma 6.6, that $\|\mathbf{M}\|_{\mathcal{S}_p}^p \ge \|\mathbf{A}\|_{\mathcal{S}_p}^p$ and thus

$$\|\mathbf{A}u\|_{2}^{p} \geq \|\mathbf{A}\|_{S_{p}}^{p} - \left(1 + O(\epsilon^{2p/3})\right) \|\mathbf{A} - \mathbf{A}_{1}\|_{S_{p}}^{p}$$

$$= \|\mathbf{A}\|_{op}^{p} - O(\epsilon^{2p/3}) \|\mathbf{A} - \mathbf{A}_{1}\|_{S_{p}}^{p}$$

$$\geq \|\mathbf{A}\|_{op}^{p} - O(\epsilon \|\mathbf{A} - \mathbf{A}_{1}\|_{S_{p}}^{p})$$

$$\geq \|\mathbf{A}\|_{op}^{p} - O(\epsilon^{2/3})$$
(6.20)

where the second to last inequality follows from recalling $p \ge 2$. The remainder of the proof is as in that following (6.7) in the proof of Theorem 6.1.

Acknowledgments: A. Bakshi and D. Woodruff would like to thank the National Science Foundation under Grant No. CCF-1815840, Office of Naval Research (ONR) grant N00014-18-1-2562, and a Simons Investigator Award. Part of this work was done while A. Bakshi was an intern at IBM Almaden. The authors thank Praneeth Kacham for pointing out an error in a previous version and for suggesting a shorter proof of Lemma 5.9. The authors also thank anonymous reviewers for their careful reading of our manuscript and for several insightful suggestions.

References

- [ACW17] Haim Avron, Kenneth L. Clarkson, and David P. Woodruff. Sharper bounds for regularized data fitting, 2017.
- [AK12] Koenraad MR Audenaert and Fuad Kittaneh. Problems and conjectures in matrix and operator inequalities. *arXiv preprint arXiv:1201.5232*, 2012.
- [AN13] Alexandr Andoni and Huy L. Nguyen. Eigenvalues of a matrix in the streaming model. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1729–1737. Society for Industrial and Applied Mathematics, 2013.
- [Ara90] Huzihiro Araki. On an inequality of Lieb and Thirring. LMaPh, 19(2):167–170, 1990.
- [Aud08] Koenraad MR Audenaert. On a norm compression inequality for 2× N partitioned block matrices. *Linear algebra and its applications*, 428(4):781–795, 2008.

- [Avr10] Haim Avron. Counting triangles in large graphs using randomized matrix trace estimation. In *Workshop on Large-scale Data Mining: Theory and Applications*, volume 10, pages 10–9, 2010.
- [BBB⁺19] Frank Ban, Vijay Bhattiprolu, Karl Bringmann, Pavel Kolev, Euiwoong Lee, and David P Woodruff. A PTAS for lp-low rank approximation. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 747–766. SIAM, 2019.
- [BBK⁺21] Ainesh Bakshi, Chiranjib Bhattacharyya, Ravi Kannan, David P Woodruff, and Samson Zhou. Learning a latent simplex in input-sparsity time. *arXiv preprint* arXiv:2105.08005, 2021.
- [BCW20] Ainesh Bakshi, Nadiia Chepurko, and David P Woodruff. Robust and sample optimal algorithms for PSD low rank approximation. In 2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS), pages 506–516. IEEE, 2020.
- [BDN15] Jean Bourgain, Sjoerd Dirksen, and Jelani Nelson. Toward a unified theory of sparse dimensionality reduction in euclidean space. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing, STOC 2015, Portland, OR, USA, June 14-17, 2015*, pages 499–508, 2015.
- [BFG96] Zhaojun Bai, Gark Fahey, and Gene Golub. Some large-scale matrix computation problems. *Journal of Computational and Applied Mathematics*, 74(1-2):71–89, 1996.
- [Bha13] Rajendra Bhatia. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.
- [BHSW20] Mark Braverman, Elad Hazan, Max Simchowitz, and Blake Woodworth. The gradient complexity of linear regression. In *Conference on Learning Theory*, pages 627–647. PMLR, 2020.
- [BKKS19] Vladimir Braverman, Robert Krauthgamer, Aditya Krishnan, and Roi Sinoff. Schatten norms in matrix streams: Hello sparsity, goodbye dimension. *arXiv preprint* arXiv:1907.05457, 2019.
- [BKL02] Rajendra Bhatia, William Kahan, and Ren-Cang Li. Pinchings and norms of scaled triangular matrices. *Linear and Multilinear Algebra*, 50(1):15–21, 2002.
- [BW18] Ainesh Bakshi and David Woodruff. Sublinear time low-rank approximation of distance matrices. In *Advances in Neural Information Processing Systems*, pages 3782–3792, 2018.
- [BWZ16] Christos Boutsidis, David P Woodruff, and Peilin Zhong. Optimal principal component analysis in distributed and streaming models. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 236–249, 2016.

- [BWZ19] Frank Ban, David Woodruff, and Qiuyi Zhang. Regularized weighted low rank approximation. *arXiv preprint arXiv:1911.06958*, 2019.
- [CCHW20] Nadiia Chepurko, Kenneth L Clarkson, Lior Horesh, and David P Woodruff. Quantum-inspired algorithms from randomized numerical linear algebra. *arXiv* preprint arXiv:2011.04125, 2020.
- [CLMW11] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):1–37, 2011.
- [CLW18] Nai-Hui Chia, Han-Hsuan Lin, and Chunhao Wang. Quantum-inspired sub-linear classical algorithms for solving low-rank linear systems. *arXiv preprint arXiv:1811.04852*, 2018.
- [Coh16] Michael B Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings of the twenty-seventh annual ACM-SIAM symposium on Discrete algorithms*, pages 278–287. SIAM, 2016.
- [CP10] Emmanuel J Candes and Yaniv Plan. Matrix completion with noise. *Proceedings of the IEEE*, 98(6):925–936, 2010.
- [CR09] Emmanuel J Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Foundations of Computational mathematics*, 9(6):717–772, 2009.
- [CW09] Kenneth L Clarkson and David P Woodruff. Numerical linear algebra in the streaming model. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*, pages 205–214. ACM, 2009.
- [CW13] Kenneth L Clarkson and David P Woodruff. Low rank approximation and regression in input sparsity time. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 81–90. ACM, 2013.
- [GKX19] Behrooz Ghorbani, Shankar Krishnan, and Ying Xiao. An investigation into neural net optimization via Hessian eigenvalue density. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2232–2241. PMLR, 09–15 Jun 2019.
- [GLT18] András Gilyén, Seth Lloyd, and Ewin Tang. Quantum-inspired low-rank stochastic regression with logarithmic dependence on the dimension. *arXiv preprint arXiv:1811.04909*, 2018.
- [GSLW19] András Gilyén, Yuan Su, Guang Hao Low, and Nathan Wiebe. Quantum singular value transformation and beyond: exponential improvements for quantum matrix arithmetics. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 193–204. ACM, 2019.

- [GXM⁺17] Shuhang Gu, Qi Xie, Deyu Meng, Wangmeng Zuo, Xiangchu Feng, and Lei Zhang. Weighted nuclear norm minimization and its applications to low level vision. *International journal of computer vision*, 121(2):183–208, 2017.
- [GZZF14] Shuhang Gu, Lei Zhang, Wangmeng Zuo, and Xiangchu Feng. Weighted nuclear norm minimization with application to image denoising. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2862–2869, 2014.
- [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [IVWW19] Piotr Indyk, Ali Vakilian, Tal Wagner, and David Woodruff. Sample-optimal low-rank approximation of distance matrices. *arXiv preprint arXiv:1906.00339*, 2019.
- [KH⁺09] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [KP16] Iordanis Kerenidis and Anupam Prakash. Quantum recommendation systems. *arXiv* preprint arXiv:1603.08675, 2016.
- [KV09] Ravi Kannan and Santosh S. Vempala. Spectral algorithms. *Found. Trends Theor. Comput. Sci.*, 4(3-4):157–288, 2009.
- [LC15] Sergey Loyka and Charalambos D. Charalambous. Novel matrix singular value inequalities and their applications to uncertain MIMO channels. *IEEE Trans. Inf. Theory*, 61(12):6623–6634, 2015.
- [LNW14a] Yi Li, Huy L Nguyen, and David P Woodruff. On sketching matrix norms and the top singular vector. In *Proceedings of the twenty-fifth annual ACM-SIAM symposium on Discrete algorithms*, pages 1562–1581. SIAM, 2014.
- [LNW14b] Yi Li, Huy L. Nguyen, and David P. Woodruff. Turnstile streaming algorithms might as well be linear sketches. In *Symposium on Theory of Computing, STOC 2014, New York, NY, USA, May 31 June 03, 2014*, pages 174–183, 2014.
- [LS13] Jörg Liesen and Zdenek Strakos. *Krylov subspace methods: principles and analysis*. Oxford University Press, 2013.
- [LW16a] Yi Li and David P Woodruff. On approximating functions of the singular values in a stream. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 726–739, 2016.
- [LW16b] Yi Li and David P Woodruff. Tight bounds for sketching the operator norm, Schatten norms, and subspace embeddings. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms and Techniques (APPROX/RANDOM 2016)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2016.

- [LW17] Yi Li and David P Woodruff. Embeddings of Schatten norms with applications to data streams. In 44th International Colloquium on Automata, Languages, and Programming (ICALP 2017). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2017.
- [LW20] Yi Li and David P. Woodruff. Input-sparsity low rank approximation in Schatten norm. *CoRR*, abs/2004.12646, 2020.
- [Mah90] Philip J Maher. Some operator inequalities concerning generalized inverses. *Illinois Journal of Mathematics*, 34(3):503–514, 1990.
- [Mah11] Michael W. Mahoney. Randomized algorithms for matrices and data. *Found. Trends Mach. Learn.*, 3(2):123–224, 2011.
- [MH02] John C Mason and David C Handscomb. *Chebyshev polynomials*. CRC press, 2002.
- [MM13a] Xiangrui Meng and Michael W Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 91–100. ACM, 2013.
- [MM13b] Xiangrui Meng and Michael W. Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Symposium on Theory of Computing Conference, STOC'13, Palo Alto, CA, USA, June 1-4, 2013,* pages 91–100, 2013.
- [MM15] Cameron Musco and Christopher Musco. Randomized block Krylov methods for stronger and faster approximate singular value decomposition. In *Advances in Neural Information Processing Systems*, pages 1396–1404, 2015.
- [MMMW21] Raphael A. Meyer, Cameron Musco, Christopher Musco, and David P. Woodruff. Hutch++: Optimal stochastic trace estimation. In 4th Symposium on Simplicity in Algorithms, SOSA 2021, Virtual Conference, January 11-12, 2021, pages 142–155, 2021.
- [Mut05] S. Muthukrishnan. Data streams: Algorithms and applications. *Found. Trends Theor. Comput. Sci.*, 1(2), 2005.
- [MW17a] Cameron Musco and David Woodruff. Is input sparsity time possible for kernel low-rank approximation? *Advances in Neural Information Processing Systems*, 30:4435–4445, 2017.
- [MW17b] Cameron Musco and David P. Woodruff. Sublinear time low-rank approximation of positive semidefinite matrices. In 58th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2017, Berkeley, CA, USA, October 15-17, 2017, pages 672–683, 2017.
- [MW21] Arvind V Mahankali and David P Woodruff. Optimal L1 column subset selection and a fast PTAS for low rank approximation. In *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 560–578. SIAM, 2021.

- [Nel11] Jelani Jelani Osei Nelson. *Sketching and streaming high-dimensional vectors*. PhD thesis, Massachusetts Institute of Technology, 2011.
- [NN13] Jelani Nelson and Huy L. Nguyen. OSNAP: faster numerical linear algebra algorithms via sparser subspace embeddings. In *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2013, 26-29 October, 2013, Berkeley, CA, USA,* pages 117–126, 2013.
- [Pea94] Barak A. Pearlmutter. Fast exact multiplication by the Hessian. *Neural Computation*, 6:147–160, 1994.
- [PPZ⁺20] William Peebles, John Peebles, Jun-Yan Zhu, Alexei A. Efros, and Antonio Torralba. The Hessian penalty: A weak prior for unsupervised disentanglement. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2020.
- [RFP10] Benjamin Recht, Maryam Fazel, and Pablo A Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM review*, 52(3):471–501, 2010.
- [Riv20] Theodore J Rivlin. *Chebyshev polynomials*. Courier Dover Publications, 2020.
- [RS00] Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- [RSML18] Patrick Rebentrost, Adrian Steffens, Iman Marvian, and Seth Lloyd. Quantum singular-value decomposition of nonsparse low-rank matrices. *Physical review A*, 97(1):012327, 2018.
- [RSW16] Ilya Razenshteyn, Zhao Song, and David P Woodruff. Weighted low rank approximations with provable guarantees. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 250–263, 2016.
- [RWYZ21] Cyrus Rashtchian, David P. Woodruff, Peng Ye, and Hanlin Zhu. Average-case communication complexity of statistical problems, 2021.
- [RWZ20] Cyrus Rashtchian, David P. Woodruff, and Hanlin Zhu. Vector-matrix-vector queries for solving linear algebra, statistics, and graph problems. In *Approximation*, *Randomization*, and Combinatorial Optimization. Algorithms and Techniques, APPROX/RANDOM 2020, August 17-19, 2020, Virtual Conference, pages 26:1–26:20, 2020.
- [Saa81] Yousef Saad. Krylov subspace methods for solving large unsymmetric linear systems. *Mathematics of computation*, 37(155):105–126, 1981.
- [SAR18] Max Simchowitz, Ahmed El Alaoui, and Benjamin Recht. Tight query complexity lower bounds for PCA via finite sample deformed wigner law. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2018, Los Angeles, CA, USA, June 25-29, 2018,* pages 1249–1259, 2018.

- [Sch60] Robert Schatten. Norm ideals of completely continuous operators. 1960.
- [SJ03] Nathan Srebro and Tommi Jaakkola. Weighted low-rank approximations. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pages 720–727, 2003.
- [SW19] Xiaofei Shi and David P. Woodruff. Sublinear time numerical linear algebra for structured matrices. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 February 1, 2019.*, pages 4918–4925, 2019.
- [SWYZ19] Xiaoming Sun, David P. Woodruff, Guang Yang, and Jialin Zhang. Querying a matrix through matrix-vector products. In 46th International Colloquium on Automata, Languages, and Programming, ICALP 2019, July 9-12, 2019, Patras, Greece, pages 94:1–94:16, 2019.
- [SWZ17] Zhao Song, David P Woodruff, and Peilin Zhong. Low rank approximation with entrywise l1-norm error. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing*, pages 688–701, 2017.
- [SWZ20] Zhao Song, David P Woodruff, and Peilin Zhong. Average case column subset selection for entrywise l1-norm loss. *arXiv preprint arXiv:2004.07986*, 2020.
- [Tan19] Ewin Tang. A quantum-inspired classical algorithm for recommendation systems. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 217–228. ACM, 2019.
- [Tao12] Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Soc., 2012.
- [Tao20] Terence Tao. Notes 3a: Eigenvalues and sums of hermitian matrices, 2020.
- [TDSL00] Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- [Tso08] Charalampos E Tsourakakis. Fast counting of triangles in large real networks without counting: Algorithms and laws. In 2008 Eighth IEEE International Conference on Data Mining, pages 608–617. IEEE, 2008.
- [VN37] John Von Neumann. Some matrix-inequalities and metrization of matric space. 1937.
- [Woo14] David P. Woodruff. Sketching as a tool for numerical linear algebra. *Foundations and Trends*® *in Theoretical Computer Science*, 10(1–2):1–157, 2014.

- [WWZ14] Karl Wimmer, Yi Wu, and Peng Zhang. Optimal query complexity for estimating the trace of a matrix. In *Automata, Languages, and Programming 41st International Colloquium, ICALP 2014, Copenhagen, Denmark, July 8-11, 2014, Proceedings, Part I*, pages 1051–1062, 2014.
- [XCS10] Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust PCA via outlier pursuit. *arXiv preprint arXiv:1010.4237*, 2010.
- [YGKM20] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael Mahoney. PyHessian: Neural networks through the lens of the Hessian, 2020.
- [YPCC16] Xinyang Yi, Dohyung Park, Yudong Chen, and Constantine Caramanis. Fast algorithms for robust PCA via gradient descent. In *Advances in neural information processing systems*, pages 4152–4160, 2016.
- [YZ16] Ming Yuan and Cun-Hui Zhang. On tensor completion via nuclear norm minimization. *Foundations of Computational Mathematics*, 16(4):1031–1068, 2016.

A Extending Prior Work on Lower Bounds

In this section, we briefly discuss prior work on estimating top singular/eigenvalues in the matrix-vector product model and why existing approaches do not immediately imply a lower bound for low-rank approximation, under any unitarily invariant norm, including Frobenius and spectral norm.

In a sequence of works, Braverman, Hazan, Simchowitz and Woodworth [BHSW20] and Simchowitz, Alaoui and Recht [SAR18] establish eigenvalue estimation lower bounds in the matrix-vector query model. We draw on their techniques and use the hard instance at the heart of their lower bound, but require additional techniques to obtain a lower bound for low-rank approximation.

The main theorem (Theorem 2.2 of [SAR18]), for k = 1, states that any randomized algorithm which outputs a vector v such that with constant probability

$$v^{\top} |\mathbf{A}| v >= (1 - O(\text{gap})) ||\mathbf{A}||_{\text{op}},$$

requires $\Omega\left(1/\sqrt{\text{gap}}\right)$ matrix-vector products, where $|\mathbf{A}|=(\mathbf{A}^2)^{1/2}$ has the same singular values as \mathbf{A} and gap $\in (0,1)$. However, this guarantee is too weak to imply a lower bound for spectral low-rank approximation.

Indeed, for this theorem to be meaningful in our setting, we require setting gap = $\Theta(\epsilon)$. However, there exist input matrices \mathbf{A} , e.g., $\mathbf{A} = \operatorname{diag}(1 + \epsilon, 1, \dots, 1, 0)$, and vector $v = \Theta\left(\sqrt{\epsilon}\right)e_1 + \left((1 - \Theta(\epsilon))e_n \right)$ such that

$$\|\mathbf{A}(\mathbf{I} - vv^{\top})\|_{\text{op}} \le (1 + \epsilon) \, \sigma_2(\mathbf{A}),$$

i.e. v yields a valid low-rank approximation but $v^{\top} \mathbf{A} v$ is only $\Theta(\epsilon)$. Note, here the gap is $\Theta(1)$ instead of the required $1 - \epsilon$ and thus we obtain no lower bound for spectral low-rank approximation.

Moreover, it can be shown that when **A** is the hard instance considered in [SAR18], i.e. **A** = $\mathbf{G} + \lambda u u^T$, where **G** is a Gaussian Orthogonal Ensemble (GOE) and u is a random unit vector on the sphere, there exists a vector v that does not satisfy the guarantee of Theorem 2.2, yet yields a spectral low-rank approximation. In particular, consider $v = \Theta(\sqrt{\epsilon})r_1 + (1 - \Theta(\epsilon))r_d$ where r_1 is the largest singular vector of $|\mathbf{A}|$ and r_d is the smallest singular vector. Since the smallest O(1) singular values of a $d \times d$ GOE can be shown to be O(1/d), and **A** is a rank-1 perturbation of a GOE, similar to the diagonal case above, we can show

$$\|\mathbf{A}(\mathbf{I} - vv^T)\|_{op} \le (1 + \epsilon) \sigma_2(\mathbf{A}),$$

yet $v^{\top} |\mathbf{A}| v$ is only $\Theta(\epsilon)$. Therefore, it is not possible to obtain a lower bound for low-rank approximation from Theorem 2.2 in a black-box manner.

B Low Rank Approximation of Matrix Polynomials

We note that polynomials of matrices are implicitly defined, even in the RAM model, and computing them explicitly would be prohibitively expensive and may destroy any sparsity structure. The proof just follows from running our algorithm on $\mathbf{M} = (\mathbf{A}^{\top}\mathbf{A})^{\ell}$. It is straightforward to simulate a matrix-vector product of the form $\mathbf{M}v$ using access to matrix-vector products for \mathbf{A} and \mathbf{A}^{\top} with an $O(\ell)$ overhead.

Theorem B.1 (Low Rank Approximation of Matrix Polynomials). Given an $n \times d$ matrix \mathbf{A} , $\ell \in \mathbb{N}$, target rank k and an accuracy parameter $\varepsilon > 0$, let $\mathbf{M} = (\mathbf{A}^{\top}\mathbf{A})^{\ell}$ or $\mathbf{M} = \mathbf{A}(\mathbf{A}^{\top}\mathbf{A})^{\ell}$. Then, for any $p \ge 1$, there exists an algorithm that uses at most $O(k\ell \log(nk)p^{1/6}/\varepsilon^{1/3})$ matrix-vector products and with probability at least 9/10 outputs a matrix $\mathbf{Z} \in \mathbb{R}^{d \times k}$ with orthonormal columns such that,

$$\left\| \mathbf{M} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p} \leq \left(1 + \varepsilon \right) \min_{\mathbf{U} : \ \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k} \left\| \mathbf{M} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{\mathcal{S}_p}.$$

The only prior work we are aware of is the algorithm of [MM15], which would achieve a worse $O(k\ell \log(nk)/\varepsilon^{1/2})$ number of matrix-vector products for the Frobenius norm and match our guarantee for the spectral norm.

C Improved Streaming Bounds

In the streaming model, the input matrix is initialized to all zeros, and at each time step, the (i, j)-th entry is updated. The updates can be positive or negative, and the goal is to output a low-rank approximation, without storing the whole matrix. The number of passes required by our algorithm is proportional to the number of *adaptive* matrix-vector queries we require. As an immediate corollary of this observation, we obtain the following formal guarantee:

Corollary C.1 (Schatten LRA in a Stream). Given a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, a target rank $k \in [d]$, an accuracy parameter $\epsilon \in (0,1)$ and any $p \geq 1$, there exists a streaming algorithm that makes $O(\log(d/\epsilon)p^{1/6}/\epsilon^{1/3})$

passes over the input, requires $O(nk/\epsilon^{1/3})$ space, and outputs a $d \times k$ matrix **Z** with orthonormal columns such that with probability at least 9/10,

$$\left\| \mathbf{A} \left(\mathbf{I} - \mathbf{Z} \mathbf{Z}^{\top} \right) \right\|_{\mathcal{S}_p}^p \leq \left(1 + \epsilon \right) \min_{\mathbf{U}: \ \mathbf{U}^{\top} \mathbf{U} = \mathbf{I}_k} \left\| \mathbf{A} \left(\mathbf{I} - \mathbf{U} \mathbf{U}^{\top} \right) \right\|_{\mathcal{S}_p}^p.$$

The only prior work on low-rank approximation in a stream is by Boutsidis, Woodruff and Zhong, who consider the special case of p=2 [BWZ16]. They obtain a single pass algorithm that requires $O(nk/\epsilon + \text{poly}(k/\epsilon))$ space and a two pass algorithm that requires $O(nk + \text{poly}(k/\epsilon))$ space. For general p, we note that recent work by Li and Woodruff [LW20] can be used to derive a streaming algorithm that obtains a worse space dependence but only requires a single pass: for $1 \le p < 2$, the space required is $\tilde{O}\left(n\left(\frac{k+k^{2/p}}{\epsilon^2} + \frac{k^{2/p}}{\epsilon^{1+2/p}}\right)\right)$ and for p > 2, the space required is $\tilde{O}\left(n\left(\frac{kn^{1-2/p}}{\epsilon^2} + \frac{k^{2/p}+n^{1-2/p}}{\epsilon^{2+2/p}}\right)\right)$.

We note that for p < 2, we obtain a polynomially better dependence on ϵ and for p > 2, the space complexity of our algorithm is linear in n, as compared to $n^{2-2/p}$ above. The optimal space complexity of Schatten-p low-rank approximation (for $p \ne 2$) in a single pass remains open.